# Decision Making Under Uncertainty: Lecture 2—Sample Average Approximation

Lecture 2
Ryan Cory-Wright
Spring 2024

## Some Housekeeping

- Reminder: Please name paper you are presenting for critical paper review and week you are presenting in (by email to me) by Friday.

## Some Housekeeping

- Reminder: Please name paper you are presenting for critical paper review and week you are presenting in (by email to me) by Friday.
- HW1 is now out, due in 2 weeks (see Insendi)—brief discussion of HW questions.

## Some Housekeeping

- Reminder: Please name paper you are presenting for critical paper review and week you are presenting in (by email to me) by Friday.
- HW1 is now out, due in 2 weeks (see Insendi)—brief discussion of HW questions.
  - Please use office hours, and don't leave it to the last minute.

## Some Housekeeping

- Reminder: Please name paper you are presenting for critical paper review and week you are presenting in (by email to me) by Friday.
- HW1 is now out, due in 2 weeks (see Insendi)—brief discussion of HW questions.
  - Please use office hours, and don't leave it to the last minute.
- My next office hours are today 3-4pm.

## Warm-Up: Solve This Problem

$$\min_{x_1, x_2} \quad x_1 + x_2$$

$$\text{s.t.} \quad \omega_1 x_1 + x_2 \geq 7$$

$$\omega_2 x_1 + x_2 \geq 4$$

$$x_1, x_2 \geq 0$$

Where $\omega_1 \sim \mathcal{U}[1, 4], \omega_2 \sim \mathcal{U}[1/3, 1]$

$$\min_{x_1, x_2} \quad x_1 + x_2$$

$$\text{s.t.} \quad \omega_1 x_1 + x_2 \geq 7$$

$$\omega_2 x_1 + x_2 \geq 4$$

$$x_1, x_2 \geq 0$$

Where $\omega_1 \sim \mathcal{U}[1, 4], \omega_2 \sim \mathcal{U}[1/3, 1]$

This problem is not well-enough defined to solve

$$\min_{x_1, x_2} \quad x_1 + x_2$$

$$\text{s.t.} \quad \omega_1 x_1 + x_2 \geq 7$$

$$\omega_2 x_1 + x_2 \geq 4$$

$$x_1, x_2 \geq 0$$

Where $\omega_1 \sim \mathcal{U}[1, 4], \omega_2 \sim \mathcal{U}[1/3, 1]$

This problem is not well-enough defined to solve

First, we don't know how $\omega_1, \omega_2$ depend on each other.

## Warm-Up: Solve This Problem

$$\min_{x_1,x_2} \quad x_1 + x_2$$

$$\text{s.t.} \quad \omega_1 x_1 + x_2 \geq 7$$

$$\omega_2 x_1 + x_2 \geq 4$$

$$x_1, x_2 \geq 0$$

Where $\omega_1 \sim \mathcal{U}[1,4], \omega_2 \sim \mathcal{U}[1/3,1]$ are indept

This problem is not well-enough defined to solve

First, we don't know how $\omega_1, \omega_2$ depend on each other. Assume indept

$$\min_{x_1, x_2} \quad x_1 + x_2$$

$$\text{s.t.} \quad \omega_1 x_1 + x_2 \geq 7, \ \omega_2 x_1 + x_2 \geq 4, \ x_1, x_2 \geq 0$$

Where $\omega_1 \sim \mathcal{U}[1, 4], \omega_2 \sim \mathcal{U}[1/3, 1]$ are indept

This problem is not well-enough defined to solve

First, we don't know how $\omega_1, \omega_2$ depend on each other. Assume indept
Second, we don't know how $x_1, x_2$ depend on $\omega$:

$$\min_{x_1, x_2} \quad x_1 + x_2$$

$$\text{s.t.} \quad \omega_1 x_1 + x_2 \geq 7, \ \omega_2 x_1 + x_2 \geq 4, \ x_1, x_2 \geq 0$$

Where $\omega_1 \sim \mathcal{U}[1, 4], \omega_2 \sim \mathcal{U}[1/3, 1]$ are indept

This problem is not well-enough defined to solve

First, we don't know how $\omega_1, \omega_2$ depend on each other. Assume indept
Second, we don't know how $x_1, x_2$ depend on $\omega$:

- Do we pick $x$, then Nature picks $\omega$, or vice versa?

$$\min_{x_1, x_2} \quad x_1 + x_2$$

$$\text{s.t.} \quad \omega_1 x_1 + x_2 \geq 7, \ \omega_2 x_1 + x_2 \geq 4, \ x_1, x_2 \geq 0$$

Where $\omega_1 \sim \mathcal{U}[1, 4], \omega_2 \sim \mathcal{U}[1/3, 1]$ are indept

This problem is not well-enough defined to solve

First, we don't know how $\omega_1$, $\omega_2$ depend on each other. Assume indept
Second, we don't know how $x_1, x_2$ depend on $\omega$:

- Do we pick $x$, then Nature picks $\omega$, or vice versa?
- First case: want to be feasible w.p.1., so minimizing $x_1 + x_2$ with $x_1 + x_2 \geq 7$, giving optimal solution of $(0, 7)$ with cost 7

$$\min_{x_1, x_2} \quad x_1 + x_2$$

$$\text{s.t.} \quad \omega_1 x_1 + x_2 \geq 7, \ \omega_2 x_1 + x_2 \geq 4, \ x_1, x_2 \geq 0$$

Where $\omega_1 \sim \mathcal{U}[1, 4], \omega_2 \sim \mathcal{U}[1/3, 1]$ are indept

This problem is not well-enough defined to solve

First, we don't know how $\omega_1$, $\omega_2$ depend on each other. Assume indept
Second, we don't know how $x_1, x_2$ depend on $\omega$:

- Do we pick $x$, then Nature picks $\omega$, or vice versa?
- First case: want to be feasible w.p.1., so minimizing $x_1 + x_2$ with $x_1 + x_2 \geq 7$, giving optimal solution of $(0, 7)$ with cost 7
- Second case: more complicated casewise analysis (exercise)

# Warm-Up: Solve This Problem

$$\min_{x_1, x_2} \quad x_1 + x_2$$

$$\text{s.t.} \quad \omega_1 x_1 + x_2 \geq 7, \ \omega_2 x_1 + x_2 \geq 4, \ x_1, x_2 \geq 0$$

Where $\omega_1 \sim \mathcal{U}[1, 4], \omega_2 \sim \mathcal{U}[1/3, 1]$ are indept
This problem is not well-enough defined to solve

First, we don't know how $\omega_1$, $\omega_2$ depend on each other. Assume indept
Second, we don't know how $x_1, x_2$ depend on $\omega$:

- Do we pick $x$, then Nature picks $\omega$, or vice versa?
- First case: want to be feasible w.p.1., so minimizing $x_1 + x_2$ with $x_1 + x_2 \geq 7$, giving optimal solution of $(0, 7)$ with cost 7
- Second case: more complicated casewise analysis (exercise)

Conclusion: Terminology matters, should define everything carefully!

## Outline of Lecture 2

Motivation: Ordinary Least Squares Regression

Sample Average Approximation: Theory

    Newsvendor: A Special Case That We Can Solve

    The General Problem

Sample Average Approximation: Algorithmics

Can we do Better? Ridge Regression and Sample-Average Approximation

Activities for if we Finish Early

    Suggested Readings

# Motivation: Ordinary Least Squares Regression

## Linear Regression Setup—Rearranging

Linear regression: $n$ i.i.d. observations of $p$-dimensional input vector $\boldsymbol{x}$ and output $y$, $\{(\boldsymbol{x}_i, y_i)\}_{i=1}^n$. We believe input-output follows model $y = \boldsymbol{x}^\top \boldsymbol{\beta}_{\text{true}} + \epsilon$, where $\boldsymbol{\beta}_{\text{true}}$ fixed vector, $\epsilon$ i.i.d. zero-mean noise.

## Linear Regression Setup—Rearranging

Linear regression: $n$ i.i.d. observations of $p$-dimensional input vector $\boldsymbol{x}$ and output $y$, $\{(\boldsymbol{x}_i, y_i)\}_{i=1}^n$. We believe input-output follows model $y = \boldsymbol{x}^\top \boldsymbol{\beta}_{\text{true}} + \epsilon$, where $\boldsymbol{\beta}_{\text{true}}$ fixed vector, $\epsilon$ i.i.d. zero-mean noise.

How to estimate $\beta$?

## Linear Regression Setup—Rearranging

Linear regression: $n$ i.i.d. observations of $p$-dimensional input vector $\boldsymbol{x}$ and output $y$, $\{(\boldsymbol{x}_i, y_i)\}_{i=1}^{n}$. We believe input-output follows model $y = \boldsymbol{x}^\top \boldsymbol{\beta}_{\text{true}} + \epsilon$, where $\boldsymbol{\beta}_{\text{true}}$ fixed vector, $\epsilon$ i.i.d. zero-mean noise.

How to estimate $\beta$? Typical answer: minimize OLS error

$$\hat{\boldsymbol{\beta}} \in \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \sum_{i=1}^{n} (y_i - \boldsymbol{x}_i^\top \boldsymbol{\beta})^2$$

## Linear Regression Setup—Rearranging

Linear regression: $n$ i.i.d. observations of $p$-dimensional input vector $\boldsymbol{x}$ and output $y$, $\{(\boldsymbol{x}_i, y_i)\}_{i=1}^n$. We believe input-output follows model $y = \boldsymbol{x}^\top \boldsymbol{\beta}_{\text{true}} + \epsilon$, where $\boldsymbol{\beta}_{\text{true}}$ fixed vector, $\epsilon$ i.i.d. zero-mean noise.

How to estimate $\beta$? Typical answer: minimize OLS error

$$\hat{\boldsymbol{\beta}} \in \arg\min_{\boldsymbol{\beta} \in \mathbb{R}^p} \sum_{i=1}^n (y_i - \boldsymbol{x}_i^\top \boldsymbol{\beta})^2$$

After some calculus

$$\hat{\boldsymbol{\beta}} = (\boldsymbol{X}^\top \boldsymbol{X})^\dagger \boldsymbol{X}^\top \boldsymbol{y},$$

where $\boldsymbol{A}^\dagger$ denotes pseudoinverse of $\boldsymbol{A}$.

## Linear Regression Setup—Rearranging

Linear regression: $n$ i.i.d. observations of $p$-dimensional input vector $\boldsymbol{x}$ and output $y$, $\{(\boldsymbol{x}_i, y_i)\}_{i=1}^n$. We believe input-output follows model $y = \boldsymbol{x}^\top \boldsymbol{\beta}_{\text{true}} + \epsilon$, where $\boldsymbol{\beta}_{\text{true}}$ fixed vector, $\epsilon$ i.i.d. zero-mean noise.

How to estimate $\beta$? Typical answer: minimize OLS error

$$\hat{\boldsymbol{\beta}} \in \arg\min_{\boldsymbol{\beta} \in \mathbb{R}^p} \sum_{i=1}^n (y_i - \boldsymbol{x}_i^\top \boldsymbol{\beta})^2$$

After some calculus

$$\hat{\boldsymbol{\beta}} = (\boldsymbol{X}^\top \boldsymbol{X})^\dagger \boldsymbol{X}^\top \boldsymbol{y},$$

where $\boldsymbol{A}^\dagger$ denotes pseudoinverse of $\boldsymbol{A}$. Assume $p$ fixed, $n > p$

$$\hat{\boldsymbol{\beta}} = (\boldsymbol{X}^\top \boldsymbol{X})^\dagger \boldsymbol{X}^\top \boldsymbol{y} \underbrace{=}_{\text{substitute } \boldsymbol{y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}} \boldsymbol{\beta}_{\text{true}} + (\boldsymbol{X}^\top \boldsymbol{X})^\dagger \boldsymbol{X}^\top \boldsymbol{\epsilon}$$

## Aside: Matrix Pseudoinverses

If $\boldsymbol{X}$ a matrix with singular value decomposition $\boldsymbol{X} = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^\top$

Then $\boldsymbol{X} = \boldsymbol{U}\boldsymbol{\Sigma}^\dagger\boldsymbol{V}^\top$ where $\boldsymbol{\Sigma}^\dagger$ is a diagonal matrix where we invert all non-zero diagonal entries, keep zeroes as zeroes.

For a symmetric matrix like $\boldsymbol{X}^\top\boldsymbol{X}$, can define

$$(\boldsymbol{X}^\top\boldsymbol{X})^\dagger := \lim_{\lambda \to 0}(\boldsymbol{X}^\top\boldsymbol{X} + \lambda\mathbb{I})^{-1}.$$

See the book "Matrix Analysis" by Horn and Johnson.

## Reminder: Almost Sure Convergence

### Almost Sure Definition

Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space and let $\{\boldsymbol{X}_i\}_{i \in \mathbb{N}}, \boldsymbol{X}$ be random variables. Suppose that $\boldsymbol{A} \in \mathcal{F}$ is a measurable set such that $\mathbb{P}(\boldsymbol{A}) = 1$ and for all $\omega \in \mathcal{A}$ we have

$$\boldsymbol{X}_i(\omega) \to \boldsymbol{X}(\omega).$$

Then, we say that $\boldsymbol{X}_i \overset{a.s.}{\to} \boldsymbol{X}$.

## Reminder: Continuous Mapping Theorem

**Continuous Mapping Theorem**

Let $\boldsymbol{X}_i, \boldsymbol{X}$ be random variables. Suppose that $\boldsymbol{X}_i \overset{a.s.}{\to} \boldsymbol{X}$ and $f$ is continuous almost everywhere. Then

$$f(\boldsymbol{X}_i) \overset{a.s.}{\to} f(\boldsymbol{X})$$

## Asymptotics of Linear Regression

Consider our rearranged equation:

$$\hat{\boldsymbol{\beta}} = (\boldsymbol{X}^\top \boldsymbol{X})^\dagger \boldsymbol{X}^\top \boldsymbol{y} = \boldsymbol{\beta}_{\text{true}} + (\boldsymbol{X}^\top \boldsymbol{X})^\dagger \boldsymbol{X}^\top \boldsymbol{\epsilon}$$

As $n \to \infty$, what happens to $\hat{\boldsymbol{\beta}}$?

## Asymptotics of Linear Regression

Consider our rearranged equation:

$$\hat{\beta} = (X^\top X)^\dagger X^\top y = \beta_{\text{true}} + (X^\top X)^\dagger X^\top \epsilon$$

As $n \to \infty$, what happens to $\hat{\beta}$?

- SLLN $\frac{1}{n} XX^\top \overset{a.s.}{\to} \mathbb{E}[x_i x_i^\top]$
- SLLN $\frac{1}{n} X^\top \epsilon \overset{a.s.}{\to} 0$

## Asymptotics of Linear Regression

Consider our rearranged equation:

$$\hat{\beta} = (\boldsymbol{X}^\top \boldsymbol{X})^\dagger \boldsymbol{X}^\top \boldsymbol{y} = \beta_{\text{true}} + (\boldsymbol{X}^\top \boldsymbol{X})^\dagger \boldsymbol{X}^\top \boldsymbol{\epsilon}$$

As $n \to \infty$, what happens to $\hat{\beta}$?

- SLLN $\frac{1}{n} \boldsymbol{X} \boldsymbol{X}^\top \overset{a.s.}{\to} \mathbb{E}[\boldsymbol{x}_i \boldsymbol{x}_i^\top]$
- SLLN $\frac{1}{n} \boldsymbol{X}^\top \boldsymbol{\epsilon} \overset{a.s.}{\to} \boldsymbol{0}$
- Therefore $\hat{\beta} \overset{a.s.}{\to} \beta_{\text{true}}$

## Asymptotics of Linear Regression

Consider our rearranged equation:

$$\hat{\beta} = (\boldsymbol{X}^\top \boldsymbol{X})^\dagger \boldsymbol{X}^\top \boldsymbol{y} = \beta_{\text{true}} + (\boldsymbol{X}^\top \boldsymbol{X})^\dagger \boldsymbol{X}^\top \epsilon$$

As $n \to \infty$, what happens to $\hat{\beta}$?

- SLLN $\frac{1}{n}\boldsymbol{X}\boldsymbol{X}^\top \overset{a.s.}{\to} \mathbb{E}[\boldsymbol{x}_i \boldsymbol{x}_i^\top]$
- SLLN $\frac{1}{n}\boldsymbol{X}^\top \epsilon \overset{a.s.}{\to} \boldsymbol{0}$
- Therefore $\hat{\beta} \overset{a.s.}{\to} \beta_{\text{true}}$



All that for a proof of convergence

**Figure 1:** Thanos Explains Empirical Risk Minimization

## What did we just do?

- We solved our first stochastic optimization problem!

## What did we just do?

- We solved our first stochastic optimization problem!
- Given sample of $n$ data points $(\boldsymbol{x}_i, y_i)$, estimate model $\boldsymbol{\beta}$ by (1) writing down stochastic optimization problem

$$\hat{\beta} = \arg \min_{\boldsymbol{\beta}} \mathbb{E}_{\boldsymbol{x},y}[(y - \boldsymbol{x}^\top \boldsymbol{\beta})^2]$$

## What did we just do?

- We solved our first stochastic optimization problem!
- Given sample of $n$ data points $(\mathbf{x}_i, y_i)$, estimate model $\boldsymbol{\beta}$ by (1) writing down stochastic optimization problem

$$\hat{\beta} = \arg\min_{\boldsymbol{\beta}} \mathbb{E}_{\mathbf{x}, y}[(y - \mathbf{x}^\top \boldsymbol{\beta})^2]$$

find estimator with least variance, (2) treating each obs. as equally likely, replacing expectation with sample-average approximation

$$\min_{\boldsymbol{\beta}} \sum_{i=1}^{n} \frac{1}{n}(y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2$$

## What did we just do?

- We solved our first stochastic optimization problem!
- Given sample of $n$ data points $(\mathbf{x}_i, y_i)$, estimate model $\boldsymbol{\beta}$ by (1) writing down stochastic optimization problem

$$\hat{\beta} = \arg\min_{\boldsymbol{\beta}} \mathbb{E}_{\mathbf{x},y}[(y - \mathbf{x}^\top \boldsymbol{\beta})^2]$$

  find estimator with least variance, (2) treating each obs. as equally likely, replacing expectation with sample-average approximation

$$\min_{\boldsymbol{\beta}} \sum_{i=1}^n \frac{1}{n}(y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2$$

- We showed $\hat{\beta}$ almost surely converges to $\beta_{\text{true}}$ as $n \to \infty$

## What did we just do?

- We solved our first stochastic optimization problem!
- Given sample of $n$ data points $(\mathbf{x}_i, y_i)$, estimate model $\boldsymbol{\beta}$ by (1) writing down stochastic optimization problem

$$\hat{\beta} = \arg\min_{\boldsymbol{\beta}} \mathbb{E}_{\mathbf{x},y}[(y - \mathbf{x}^\top \boldsymbol{\beta})^2]$$

find estimator with least variance, (2) treating each obs. as equally likely, replacing expectation with sample-average approximation

$$\min_{\boldsymbol{\beta}} \sum_{i=1}^{n} \frac{1}{n}(y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2$$

- We showed $\hat{\beta}$ almost surely converges to $\beta_{\text{true}}$ as $n \to \infty$
- So supervised learning is special case of stochastic optimization!

## What did we just do?

- We solved our first stochastic optimization problem!
- Given sample of $n$ data points $(\mathbf{x}_i, y_i)$, estimate model $\boldsymbol{\beta}$ by (1) writing down stochastic optimization problem

$$\hat{\beta} = \arg \min_{\boldsymbol{\beta}} \mathbb{E}_{\mathbf{x}, y}[(y - \mathbf{x}^\top \boldsymbol{\beta})^2]$$

find estimator with least variance, (2) treating each obs. as equally likely, replacing expectation with sample-average approximation

$$\min_{\boldsymbol{\beta}} \sum_{i=1}^{n} \frac{1}{n}(y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2$$

- We showed $\hat{\beta}$ almost surely converges to $\beta_{\text{true}}$ as $n \to \infty$
- So supervised learning is special case of stochastic optimization!
- This would take a ML class 3-4 lectures; let's take a breath here!

## What did we just do?

- We solved our first stochastic optimization problem!
- Given sample of $n$ data points $(\boldsymbol{x}_i, y_i)$, estimate model $\boldsymbol{\beta}$ by (1) writing down stochastic optimization problem

$$\hat{\beta} = \arg \min_{\boldsymbol{\beta}} \mathbb{E}_{\boldsymbol{x}, y}[(y - \boldsymbol{x}^\top \boldsymbol{\beta})^2]$$

  find estimator with least variance, (2) treating each obs. as equally likely, replacing expectation with sample-average approximation

$$\min_{\boldsymbol{\beta}} \sum_{i=1}^{n} \frac{1}{n}(y_i - \boldsymbol{x}_i^\top \boldsymbol{\beta})^2$$

- We showed $\hat{\beta}$ almost surely converges to $\beta_{\text{true}}$ as $n \to \infty$
- So supervised learning is special case of stochastic optimization!
- This would take a ML class 3-4 lectures; let's take a breath here!
- Plan for lecture: Show holds more generally, how to solve SAA

**Let's break for five minutes here**

# Sample Average Approximation: Theory

Let's warm up with a special case

## Hot off the Press: The Newsvendor Problem

- A newsvendor (newspaper salesperson) needs to decide how many newspapers $x$ to buy to maximize their profit

## Hot off the Press: The Newsvendor Problem

- A newsvendor (newspaper salesperson) needs to decide how many newspapers $x$ to buy to maximize their profit
- She doesn't know how many newspapers there are demand for, $D_\omega$ in scenario $\omega$. But she does know the probability distribution of $D_\omega$

## Hot off the Press: The Newsvendor Problem

- A newsvendor (newspaper salesperson) needs to decide how many newspapers $x$ to buy to maximize their profit
- She doesn't know how many newspapers there are demand for, $D_\omega$ in scenario $\omega$. But she does know the probability distribution of $D_\omega$
- Each newspaper costs $c$, can be sold for $q$ if there is demand
- Unsold newspapers get thrown in the recycling bin
- How to optimally set $x$?

## Hot off the Press: The Newsvendor Problem

$$\max_{x \geq 0} \mathbb{E}_{\omega}[\min(D_{\omega}, x)q - cx]$$

## Hot off the Press: The Newsvendor Problem

$$\max_{x \geq 0} \mathbb{E}_\omega[\min(D_\omega, x)q - cx]$$

Two cases: $x > D_\omega$ or $x \leq D_\omega$. Rewrite using conditional expectations

## Hot off the Press: The Newsvendor Problem

$$\max_{x \geq 0} \mathbb{E}_\omega[\min(D_\omega, x)q - cx]$$

Two cases: $x > D_\omega$ or $x \leq D_\omega$. Rewrite using conditional expectations

$$\max_{x \geq 0} \mathbb{E}_\omega[D_\omega q - cx | x \leq D_\omega]\mathbb{P}(x \leq D_\omega) + \mathbb{E}_\omega[qx - cx | x > D_\omega]\mathbb{P}(x > D_\omega)$$

## Hot off the Press: The Newsvendor Problem

$$\max_{x \geq 0} \mathbb{E}_\omega[\min(D_\omega, x)q - cx]$$

Two cases: $x > D_\omega$ or $x \leq D_\omega$. Rewrite using conditional expectations

$$\max_{x \geq 0} \mathbb{E}_\omega[D_\omega q - cx | x \leq D_\omega]\mathbb{P}(x \leq D_\omega) + \mathbb{E}_\omega[qx - cx | x > D_\omega]\mathbb{P}(x > D_\omega)$$

This is convex in $x$, so differentiate with respect to $x$, require that 0 in subgradient.

## Hot off the Press: The Newsvendor Problem

$$\max_{x \geq 0} \mathbb{E}_\omega[\min(D_\omega, x)q - cx]$$

Two cases: $x > D_\omega$ or $x \leq D_\omega$. Rewrite using conditional expectations

$$\max_{x \geq 0} \mathbb{E}_\omega[D_\omega q - cx | x \leq D_\omega]\mathbb{P}(x \leq D_\omega) + \mathbb{E}_\omega[qx - cx | x > D_\omega]\mathbb{P}(x > D_\omega)$$

This is convex in $x$, so differentiate with respect to $x$, require that 0 in subgradient.

Eventually get

$$x^\star \in F_{D_\omega}^{-1}\left(\frac{q - c}{q}\right)$$

14

## Hot off the Press: The Newsvendor Problem

$$\max_{x \geq 0} \mathbb{E}_{\omega}[\min(D_{\omega}, x)q - cx]$$

Two cases: $x > D_{\omega}$ or $x \leq D_{\omega}$. Rewrite using conditional expectations

$$\max_{x \geq 0} \mathbb{E}_{\omega}[D_{\omega}q - cx | x \leq D_{\omega}]\mathbb{P}(x \leq D_{\omega}) + \mathbb{E}_{\omega}[qx - cx | x > D_{\omega}]\mathbb{P}(x > D_{\omega})$$

This is convex in $x$, so differentiate with respect to $x$, require that 0 in subgradient.

Eventually get

$$x^{\star} \in F_{D_{\omega}}^{-1}\left(\frac{q - c}{q}\right)$$

That is, a $\frac{(q-c)}{q}$th quantile of $D_{\omega}$

Insight: setting $x$ equal to $\mathbb{E}[D_{\omega}]$ could be bad, especially if $q \gg c$

# The General Problem

## Overall Problem Setting: Two-Stage Stochastic Linear Opt

Consider stochastic optimization problem:

$$\min_{\boldsymbol{x} \in \mathbb{R}^n} \quad \boldsymbol{c}^\top \boldsymbol{x} + \mathbb{E}_\omega[h(\boldsymbol{x}, \boldsymbol{\omega})]$$
$$\text{s.t.} \quad \boldsymbol{A}\boldsymbol{x} \le \boldsymbol{b}$$

## Overall Problem Setting: Two-Stage Stochastic Linear Opt

Consider stochastic optimization problem:

$$\min_{\mathbf{x} \in \mathbb{R}^n} \quad \mathbf{c}^\top \mathbf{x} + \mathbb{E}_\omega[h(\mathbf{x}, \omega)]$$
$$\text{s.t.} \quad \mathbf{A}\mathbf{x} \leq \mathbf{b}$$

where

$$h(\mathbf{x}, \omega) := \min_{\mathbf{y}(\omega)} \mathbf{q}(\omega)^\top \mathbf{y}(\omega)$$
$$\text{s.t.} \quad \mathbf{D}(\omega)\mathbf{x} + \mathbf{F}(\omega)\mathbf{y}(\omega) \leq \mathbf{d}(\omega)$$

## Overall Problem Setting: Two-Stage Stochastic Linear Opt

Consider stochastic optimization problem:

$$\min_{\boldsymbol{x} \in \mathbb{R}^n} \quad \boldsymbol{c}^\top \boldsymbol{x} + \mathbb{E}_{\boldsymbol{\omega}}[h(\boldsymbol{x}, \boldsymbol{\omega})]$$

$$\text{s.t.} \quad \boldsymbol{A}\boldsymbol{x} \leq \boldsymbol{b}$$

where

$$h(\boldsymbol{x}, \boldsymbol{\omega}) := \min_{\boldsymbol{y}(\boldsymbol{\omega})} \boldsymbol{q}(\boldsymbol{\omega})^\top \boldsymbol{y}(\boldsymbol{\omega})$$

$$\text{s.t.} \quad \boldsymbol{D}(\omega)\boldsymbol{x} + \boldsymbol{F}(\omega)\boldsymbol{y}(\omega) \leq \boldsymbol{d}(\omega)$$

- $x$ are our first-stage (or here-and-now) decision variables, which we select before nature picks $\boldsymbol{\omega}$

## Overall Problem Setting: Two-Stage Stochastic Linear Opt

Consider stochastic optimization problem:

$$\min_{\boldsymbol{x} \in \mathbb{R}^n} \quad \boldsymbol{c}^\top \boldsymbol{x} + \mathbb{E}_{\omega}[h(\boldsymbol{x}, \boldsymbol{\omega})]$$
$$\text{s.t.} \quad \boldsymbol{A}\boldsymbol{x} \leq \boldsymbol{b}$$

where

$$h(\boldsymbol{x}, \boldsymbol{\omega}) := \min_{\boldsymbol{y}(\boldsymbol{\omega})} \boldsymbol{q}(\boldsymbol{\omega})^\top \boldsymbol{y}(\boldsymbol{\omega})$$
$$\text{s.t.} \quad \boldsymbol{D}(\omega)\boldsymbol{x} + \boldsymbol{F}(\omega)\boldsymbol{y}(\boldsymbol{\omega}) \leq \boldsymbol{d}(\omega)$$

- $x$ are our first-stage (or here-and-now) decision variables, which we select before nature picks $\boldsymbol{\omega}$
- $\boldsymbol{\omega}$ are the random variables selected by nature, according to their joint probability distribution (assumed to be known)

## Overall Problem Setting: Two-Stage Stochastic Linear Opt

Consider stochastic optimization problem:

$$\min_{\mathbf{x} \in \mathbb{R}^n} \quad \mathbf{c}^\top \mathbf{x} + \mathbb{E}_\omega[h(\mathbf{x}, \boldsymbol{\omega})]$$
$$\text{s.t.} \quad \mathbf{A}\mathbf{x} \leq \mathbf{b}$$

where

$$h(\mathbf{x}, \boldsymbol{\omega}) := \min_{\mathbf{y}(\boldsymbol{\omega})} \mathbf{q}(\boldsymbol{\omega})^\top \mathbf{y}(\boldsymbol{\omega})$$
$$\text{s.t.} \quad \mathbf{D}(\omega)\mathbf{x} + \mathbf{F}(\omega)\mathbf{y}(\omega) \leq \mathbf{d}(\omega)$$

- $x$ are our first-stage (or here-and-now) decision variables, which we select before nature picks $\boldsymbol{\omega}$
- $\boldsymbol{\omega}$ are the random variables selected by nature, according to their joint probability distribution (assumed to be known)
- $y(\boldsymbol{\omega})$ are our second-stage (or wait-and-see, or recourse) decision variables, that we are allowed to pick after nature picks $\boldsymbol{\omega}$

Consider stochastic optimization problem:

$$\min_{\boldsymbol{x} \in \mathbb{R}^n} \quad \boldsymbol{c}^\top \boldsymbol{x} + \mathbb{E}_\omega[h(\boldsymbol{x}, \boldsymbol{\omega})]$$

$$\text{s.t.} \quad \boldsymbol{A}\boldsymbol{x} \leq \boldsymbol{b}$$

where

$$h(\boldsymbol{x}, \boldsymbol{\omega}) := \min_{\boldsymbol{y}(\boldsymbol{\omega})} \boldsymbol{q}(\boldsymbol{\omega})^\top \boldsymbol{y}(\boldsymbol{\omega})$$

$$\text{s.t.} \quad \boldsymbol{D}(\omega)\boldsymbol{x} + \boldsymbol{F}(\omega)\boldsymbol{y}(\boldsymbol{\omega}) \leq \boldsymbol{d}(\boldsymbol{\omega})$$

- $x$ are our first-stage (or here-and-now) decision variables, which we select before nature picks $\boldsymbol{\omega}$
- $\boldsymbol{\omega}$ are the random variables selected by nature, according to their joint probability distribution (assumed to be known)
- $y(\boldsymbol{\omega})$ are our second-stage (or wait-and-see, or recourse) decision variables, that we are allowed to pick after nature picks $\boldsymbol{\omega}$
- A linear optimization problem with random parameters

15

## What Makes This Problem Hard?

- Complexity Theory: Solving this problem is $\#P$-hard

## What Makes This Problem Hard?

- Complexity Theory: Solving this problem is $\#P$-hard
  - See Hanasusanto et al. (Math. Prog., 2016) for a proof

## What Makes This Problem Hard?

- Complexity Theory: Solving this problem is $\#P$-hard
  - See Hanasusanto et al. (Math. Prog., 2016) for a proof
  - Who knows what this means?

## What Makes This Problem Hard?

- Complexity Theory: Solving this problem is $\#P$-hard
    - See Hanasusanto et al. (Math. Prog., 2016) for a proof
    - Who knows what this means?
    - As hard as counting number of solutions to NP-hard problem

## What Makes This Problem Hard?

- Complexity Theory: Solving this problem is $\#P$-hard
  - See Hanasusanto et al. (Math. Prog., 2016) for a proof
  - Who knows what this means?
  - As hard as counting number of solutions to NP-hard problem
  - I once heard someone say "Judging a problem by its complexity is like judging someone by the worst thing they have ever done"
    —In and of itself, $\#P$-hard doesn't mean intractable

## What Makes This Problem Hard?

- Complexity Theory: Solving this problem is $\#P$-hard
  - See Hanasusanto et al. (Math. Prog., 2016) for a proof
  - Who knows what this means?
  - As hard as counting number of solutions to NP-hard problem
  - I once heard someone say "Judging a problem by its complexity is like judging someone by the worst thing they have ever done" —In and of itself, $\#P$-hard doesn't mean intractable
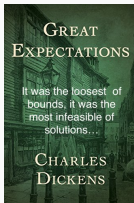- Numerically: Expectations hard to evaluate in high-dim settings



**Figure 2:** Dickens explains the curse of dimensionality

## What Makes This Problem Hard?

- Complexity Theory: Solving this problem is $\#P$-hard
  - See Hanasusanto et al. (Math. Prog., 2016) for a proof
  - Who knows what this means?
  - As hard as counting number of solutions to NP-hard problem
  - I once heard someone say "Judging a problem by its complexity is like judging someone by the worst thing they have ever done"
    —In and of itself, $\#P$-hard doesn't mean intractable
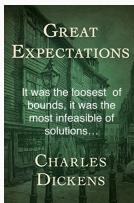- Numerically: Expectations hard to evaluate in high-dim settings



GREAT
EXPECTATIONS

It was the loosest of bounds, it was the most infeasible of solutions...

CHARLES
DICKENS

**Figure 2:** Dickens explains the curse of dimensionality

- Structure of Optimal Solutions: In general, $y$ a function of $\omega$

## Sample Average Approximation to the Rescue

Let's play same game as in the linear regression case!

## Sample Average Approximation to the Rescue

Let's play same game as in the linear regression case!
Replace (unknown) expectation over $\omega$ with expectation over empirical distribution $\omega_i$. With $n$ observations of $\omega$, or $n$ "scenarios", solve:

$$\hat{x} \in \arg\min_{x \in \mathbb{R}^n} \quad c^\top x + \frac{1}{n} \sum_{i=1}^{n} h(x, \omega^i)$$

$$\text{s.t.} \quad Ax \leq b$$

## Sample Average Approximation to the Rescue

Let's play same game as in the linear regression case!

Replace (unknown) expectation over $\omega$ with expectation over empirical distribution $\omega_i$. With $n$ observations of $\omega$, or $n$ "scenarios", solve:

$$\hat{x} \in \arg \min_{x \in \mathbb{R}^n} \quad c^\top x + \frac{1}{n} \sum_{i=1}^{n} h(x, \omega^i)$$

$$\text{s.t.} \quad Ax \leq b$$

Why is this a good thing to do?

## Sample Average Approximation to the Rescue

Let's play same game as in the linear regression case!

Replace (unknown) expectation over $\omega$ with expectation over empirical distribution $\omega_i$. With $n$ observations of $\omega$, or $n$ "scenarios", solve:

$$\hat{x} \in \arg\min_{x \in \mathbb{R}^n} \quad c^\top x + \frac{1}{n} \sum_{i=1}^n h(x, \omega^i)$$

$$\text{s.t.} \quad Ax \leq b$$

Why is this a good thing to do? Justifications:

- Joint distribution over $\omega$ only exists *in our imagination*, while empirical distribution constructed from data, which is real

### Sample Average Approximation to the Rescue

Let's play same game as in the linear regression case!
Replace (unknown) expectation over $\omega$ with expectation over empirical distribution $\omega_i$. With $n$ observations of $\omega$, or $n$ "scenarios", solve:

$$\hat{x} \in \arg \min_{x \in \mathbb{R}^n} \quad c^\top x + \frac{1}{n} \sum_{i=1}^{n} h(x, \omega^i)$$

$$\text{s.t.} \quad Ax \leq b$$

Why is this a good thing to do? Justifications:

- Joint distribution over $\omega$ only exists *in our imagination*, while empirical distribution constructed from data, which is real

- As $n \to \infty$, for i.i.d. $\omega^i$, $\hat{x}$ almost surely converges to a minimizer of our two-stage problem under true joint distribution of $\omega$

## Sample Average Approximation to the Rescue

Let's play same game as in the linear regression case!
Replace (unknown) expectation over $\omega$ with expectation over empirical distribution $\omega_i$. With $n$ observations of $\omega$, or $n$ "scenarios", solve:

$$\hat{x} \in \arg\min_{x \in \mathbb{R}^n} \quad c^\top x + \frac{1}{n} \sum_{i=1}^{n} h(x, \omega^i)$$

$$\text{s.t.} \quad Ax \leq b$$

Why is this a good thing to do? Justifications:

- Joint distribution over $\omega$ only exists *in our imagination*, while empirical distribution constructed from data, which is real

- As $n \to \infty$, for i.i.d. $\omega^i$, $\hat{x}$ almost surely converges to a minimizer of our two-stage problem under true joint distribution of $\omega$

- Who can tell me why we use "arg min" and "a minimizer" here?

## Almost Sure Convergence Proof (Sketch)

- Define a sample-average function, redefine expected value

$$\hat{g}_N(\boldsymbol{x}) := \min_{\boldsymbol{y}(\boldsymbol{\omega}^i)} \boldsymbol{c}^\top \boldsymbol{x} + \frac{1}{N} \sum_{i=1}^{n} h(\boldsymbol{x}, \boldsymbol{\omega}^i),$$

$$g(\boldsymbol{x}) := \min_{\boldsymbol{y}(\omega)} \mathbb{E}_\omega [\boldsymbol{c}^\top \boldsymbol{x} + \frac{1}{N} \sum_{i=1}^{n} h(\boldsymbol{x}, \boldsymbol{\omega})]$$

## Aside

$h$ is the optimal value of a minimization problem. Why is it convex?

## Aside

$h$ is the optimal value of a minimization problem. Why is it convex?



$$h(\boldsymbol{x}, \boldsymbol{\omega}) := \min_{\boldsymbol{y}(\omega)} \boldsymbol{q}(\boldsymbol{\omega})^\top \boldsymbol{y}(\boldsymbol{\omega}) \text{ s.t. } \boldsymbol{D}(\omega)\boldsymbol{x} + \boldsymbol{F}(\omega)\boldsymbol{y}(\boldsymbol{\omega}) \leq \boldsymbol{d}(\boldsymbol{\omega})$$

## Aside

$h$ is the optimal value of a minimization problem. Why is it convex?



$$h(\boldsymbol{x}, \boldsymbol{\omega}) := \min_{\boldsymbol{y}(\omega)} \boldsymbol{q}(\boldsymbol{\omega})^\top \boldsymbol{y}(\boldsymbol{\omega}) \text{ s.t. } \boldsymbol{D}(\omega)\boldsymbol{x} + \boldsymbol{F}(\omega)\boldsymbol{y}(\boldsymbol{\omega}) \leq \boldsymbol{d}(\boldsymbol{\omega})$$

Duality!

$$h(\boldsymbol{x}, \boldsymbol{\omega}) = \max_{\boldsymbol{\mu}(\omega)} \; (\boldsymbol{d}(\boldsymbol{\omega}) - \boldsymbol{D}(\omega)\boldsymbol{x})^\top \boldsymbol{\mu}(\boldsymbol{\omega}) \text{ s.t. } \boldsymbol{F}(\omega)^\top \boldsymbol{\mu}(\boldsymbol{\omega}) = \boldsymbol{q}(\boldsymbol{\omega}), \boldsymbol{\mu}(\boldsymbol{\omega}) \leq \boldsymbol{0}$$

## Aside

$h$ is the optimal value of a minimization problem. Why is it convex?



$$h(\boldsymbol{x}, \boldsymbol{\omega}) := \min_{\boldsymbol{y}(\boldsymbol{\omega})} \boldsymbol{q}(\boldsymbol{\omega})^\top \boldsymbol{y}(\boldsymbol{\omega}) \text{ s.t. } \boldsymbol{D}(\omega)\boldsymbol{x} + \boldsymbol{F}(\omega)\boldsymbol{y}(\boldsymbol{\omega}) \leq \boldsymbol{d}(\boldsymbol{\omega})$$

Duality!

$$h(\boldsymbol{x}, \boldsymbol{\omega}) = \max_{\boldsymbol{\mu}(\omega)} \ (\boldsymbol{d}(\boldsymbol{\omega}) - \boldsymbol{D}(\omega)\boldsymbol{x})^\top \boldsymbol{\mu}(\omega) \text{ s.t. } \boldsymbol{F}(\omega)^\top \boldsymbol{\mu}(\omega) = \boldsymbol{q}(\boldsymbol{\omega}), \boldsymbol{\mu}(\omega) \leq \boldsymbol{0}$$

$h(\boldsymbol{x}, \boldsymbol{\omega})$ is the pointwise maximum of functions linear in $\boldsymbol{x}$, hence convex

## Aside

$h$ is the optimal value of a minimization problem. Why is it convex?



$$h(\boldsymbol{x}, \boldsymbol{\omega}) := \min_{\boldsymbol{y}(\boldsymbol{\omega})} \boldsymbol{q}(\boldsymbol{\omega})^\top \boldsymbol{y}(\boldsymbol{\omega}) \text{ s.t. } \boldsymbol{D}(\omega)\boldsymbol{x} + \boldsymbol{F}(\omega)\boldsymbol{y}(\boldsymbol{\omega}) \leq \boldsymbol{d}(\boldsymbol{\omega})$$

Duality!

$$h(\boldsymbol{x}, \boldsymbol{\omega}) = \max_{\boldsymbol{\mu}(\boldsymbol{\omega})} \ (\boldsymbol{d}(\boldsymbol{\omega}) - \boldsymbol{D}(\omega)\boldsymbol{x})^\top \boldsymbol{\mu}(\boldsymbol{\omega}) \text{ s.t. } \boldsymbol{F}(\omega)^\top \boldsymbol{\mu}(\boldsymbol{\omega}) = \boldsymbol{q}(\boldsymbol{\omega}), \boldsymbol{\mu}(\boldsymbol{\omega}) \leq \boldsymbol{0}$$

$h(\boldsymbol{x}, \boldsymbol{\omega})$ is the pointwise maximum of functions linear in $\boldsymbol{x}$, hence convex
Pointwise maximum also reveals $h$ is continuous on its domain

## Almost Sure Convergence Proof (Sketch)

- Define a sample-average function, redefine expected value

$$\hat{g}_N(\boldsymbol{x}) := \boldsymbol{c}^\top \boldsymbol{x} + \frac{1}{N} \sum_{i=1}^{N} h(\boldsymbol{x}, \boldsymbol{\omega}^i),$$

$$g(\boldsymbol{x}) := \boldsymbol{c}^\top \boldsymbol{x} + \mathbb{E}_\omega[h(\boldsymbol{x}, \boldsymbol{\omega})]$$

- By SLLN, continuity of $g_N, g$: $g_N(\boldsymbol{x}) \xrightarrow{a.s.} g(\boldsymbol{x}) \ \forall \boldsymbol{x} : \boldsymbol{Ax} \leq \boldsymbol{b}$

---

[1]See Corollary 3 of "Monte Carlo Sampling Methods" by Shapiro (2003) for details.

## Almost Sure Convergence Proof (Sketch)

- Define a sample-average function, redefine expected value

$$\hat{g}_N(\boldsymbol{x}) := \boldsymbol{c}^\top \boldsymbol{x} + \frac{1}{N} \sum_{i=1}^{N} h(\boldsymbol{x}, \boldsymbol{\omega}^i),$$

$$g(\boldsymbol{x}) := \boldsymbol{c}^\top \boldsymbol{x} + \mathbb{E}_\omega[h(\boldsymbol{x}, \boldsymbol{\omega})]$$

- By SLLN, continuity of $g_N, g$: $g_N(\boldsymbol{x}) \overset{a.s.}{\to} g(\boldsymbol{x}) \; \forall \boldsymbol{x} : \boldsymbol{Ax} \leq \boldsymbol{b}$
- Therefore, (under mild conditions[1]), $\inf_{\boldsymbol{x}} g_N(\boldsymbol{x}) \overset{a.s.}{\to} \inf_{\boldsymbol{x}} g(\boldsymbol{x})$

---

[1] See Corollary 3 of "Monte Carlo Sampling Methods" by Shapiro (2003) for details.

## When Things go Wrong, as They Sometimes Will

Let's look at our sample-average approximation again:

$$\hat{\boldsymbol{x}} \in \arg \min_{\boldsymbol{x} \in \mathbb{R}^n} \quad \boldsymbol{c}^\top \boldsymbol{x} + \frac{1}{n} \sum_{i=1}^n h(\boldsymbol{x}, \boldsymbol{\omega}^i)$$

$$\text{s.t.} \quad \boldsymbol{A}\boldsymbol{x} \leq \boldsymbol{b}$$

What can go wrong?

## When Things go Wrong, as They Sometimes Will

Let's look at our sample-average approximation again:

$$\hat{x} \in \arg\min_{x \in \mathbb{R}^n} \quad c^\top x + \frac{1}{n} \sum_{i=1}^{n} h(x, \omega^i)$$

$$\text{s.t.} \quad Ax \leq b$$

What can go wrong? In practice, we have a finite number of observations. That means:

## When Things go Wrong, as They Sometimes Will

Let's look at our sample-average approximation again:

$$\hat{x} \in \arg\min_{x \in \mathbb{R}^n} \quad c^\top x + \frac{1}{n} \sum_{i=1}^{n} h(x, \omega^i)$$

$$\text{s.t.} \quad Ax \leq b$$

What can go wrong? In practice, we have a finite number of observations. That means:

- $\hat{x}$ may not be feasible for unseen $\omega$'s

## When Things go Wrong, as They Sometimes Will

Let's look at our sample-average approximation again:

$$\hat{\boldsymbol{x}} \in \arg \min_{\boldsymbol{x} \in \mathbb{R}^n} \quad \boldsymbol{c}^\top \boldsymbol{x} + \frac{1}{n} \sum_{i=1}^{n} h(\boldsymbol{x}, \boldsymbol{\omega}^i)$$

$$\text{s.t.} \quad \boldsymbol{A}\boldsymbol{x} \leq \boldsymbol{b}$$

What can go wrong? In practice, we have a finite number of observations. That means:

- $\hat{\boldsymbol{x}}$ may not be feasible for unseen $\boldsymbol{\omega}$'s
    - Can include all extreme points of joint dist of $\omega$, or if $h(\boldsymbol{x}, \boldsymbol{\omega}^i)$ is (almost surely) feasible for any $\boldsymbol{x}$—(relatively) complete recourse

### When Things go Wrong, as They Sometimes Will

Let's look at our sample-average approximation again:

$$\hat{\boldsymbol{x}} \in \arg\min_{\boldsymbol{x} \in \mathbb{R}^n} \quad \boldsymbol{c}^\top \boldsymbol{x} + \frac{1}{n} \sum_{i=1}^{n} h(\boldsymbol{x}, \boldsymbol{\omega}^i)$$

$$\text{s.t.} \quad \boldsymbol{A}\boldsymbol{x} \leq \boldsymbol{b}$$

What can go wrong? In practice, we have a finite number of observations. That means:

- $\hat{\boldsymbol{x}}$ may not be feasible for unseen $\boldsymbol{\omega}$'s
    - Can include all extreme points of joint dist of $\omega$, or if $h(\boldsymbol{x}, \boldsymbol{\omega}^i)$ is (almost surely) feasible for any $\boldsymbol{x}$—(relatively) complete recourse
- $\hat{\boldsymbol{x}}_N$ might be far from $\boldsymbol{x}^\star$, especially if $N$ small relative to dim of $\boldsymbol{x}$
    - A motivation for distributionally robust optimization—see later

**Let's break for five minutes.**
**Then talk about how to solve these problems**

# Sample Average Approximation: Algorithmics

## First Strategy: Solve the Deterministic Equivalent

We can view the sample-average approximation as one big linear optimization problem and throw it to Mosek or Gurobi

## First Strategy: Solve the Deterministic Equivalent

We can view the sample-average approximation as one big linear optimization problem and throw it to Mosek or Gurobi

- Make a copy of $y^i$ for each scenario $\omega^i$ and solve

$$\hat{x} \in \arg\min_{x \in \mathbb{R}^n} \quad c^\top x + \frac{1}{n} \sum_{i=1}^n h(x, \omega^i)$$

$$\text{s.t.} \quad Ax \leq b$$

### First Strategy: Solve the Deterministic Equivalent

We can view the sample-average approximation as one big linear
optimization problem and throw it to Mosek or Gurobi

- Make a copy of $y^i$ for each scenario $\omega^i$ and solve

$$\hat{x} \in \arg\min_{x \in \mathbb{R}^n} \quad c^\top x + \frac{1}{n} \sum_{i=1}^n h(x, \omega^i)$$

$$\text{s.t.} \quad Ax \leq b$$

- Pros: very quick to code, if it works, then we are done

### First Strategy: Solve the Deterministic Equivalent

We can view the sample-average approximation as one big linear optimization problem and throw it to Mosek or Gurobi

- Make a copy of $y^i$ for each scenario $\omega^i$ and solve

$$\hat{x} \in \arg\min_{x \in \mathbb{R}^n} \quad c^\top x + \frac{1}{n} \sum_{i=1}^{n} h(x, \omega^i)$$

$$\text{s.t.} \quad Ax \le b$$

- Pros: very quick to code, if it works, then we are done
- Good first thing to try

### First Strategy: Solve the Deterministic Equivalent

We can view the sample-average approximation as one big linear optimization problem and throw it to Mosek or Gurobi

- Make a copy of $y^i$ for each scenario $\omega^i$ and solve

$$\hat{x} \in \arg\min_{x \in \mathbb{R}^n} \quad c^\top x + \frac{1}{n} \sum_{i=1}^{n} h(x, \omega^i)$$

$$\text{s.t.} \quad Ax \leq b$$

- Pros: very quick to code, if it works, then we are done
- Good first thing to try
- Cons: this optimization problem might be big. Really big

### First Strategy: Solve the Deterministic Equivalent

We can view the sample-average approximation as one big linear optimization problem and throw it to Mosek or Gurobi

- Make a copy of $y^i$ for each scenario $\omega^i$ and solve

$$\hat{x} \in \arg\min_{x \in \mathbb{R}^n} \quad c^\top x + \frac{1}{n} \sum_{i=1}^{n} h(x, \omega^i)$$

$$\text{s.t.} \quad Ax \leq b$$

- Pros: very quick to code, if it works, then we are done
- Good first thing to try
- Cons: this optimization problem might be big. Really big
- Example: electricity market with random demand at 20 nodes that can independently be "low" or "high"

### First Strategy: Solve the Deterministic Equivalent

We can view the sample-average approximation as one big linear optimization problem and throw it to Mosek or Gurobi

- Make a copy of $y^i$ for each scenario $\omega^i$ and solve

$$\hat{x} \in \arg \min_{x \in \mathbb{R}^n} \quad c^\top x + \frac{1}{n} \sum_{i=1}^{n} h(x, \omega^i)$$

$$\text{s.t.} \quad Ax \leq b$$

- Pros: very quick to code, if it works, then we are done
- Good first thing to try
- Cons: this optimization problem might be big. Really big
- Example: electricity market with random demand at 20 nodes that can independently be "low" or "high" That's $2^{20} = 1048576$ copies of $y$, which is intractable for a real market

### First Strategy: Solve the Deterministic Equivalent

We can view the sample-average approximation as one big linear optimization problem and throw it to Mosek or Gurobi

- Make a copy of $y^i$ for each scenario $\omega^i$ and solve

$$\hat{x} \in \arg\min_{x \in \mathbb{R}^n} \quad c^\top x + \frac{1}{n} \sum_{i=1}^{n} h(x, \omega^i)$$

$$\text{s.t.} \quad Ax \leq b$$

- Pros: very quick to code, if it works, then we are done
- Good first thing to try
- Cons: this optimization problem might be big. Really big
- Example: electricity market with random demand at 20 nodes that can independently be "low" or "high" That's $2^{20} = 1048576$ copies of $y$, which is intractable for a real market
- Still, you can sometimes do well by subsampling the scenarios (Shapiro and Homem-de-Mello, 1998)

## Second Strategy: Decompose the Problem

What optimizers usually do: use a decomposition scheme called Benders decomposition (sometimes called the "L-shaped" method)

## Second Strategy: Decompose the Problem

What optimizers usually do: use a decomposition scheme called Benders decomposition (sometimes called the "L-shaped" method)

Consider

$$\min_{\boldsymbol{x} \in \mathbb{R}^n} \quad \boldsymbol{c}^\top \boldsymbol{x} + \frac{1}{n} \sum_{i=1}^{n} h(\boldsymbol{x}, \boldsymbol{\omega}^i)$$

$$\text{s.t.} \quad \boldsymbol{A}\boldsymbol{x} \leq \boldsymbol{b}$$

Let $\theta \geq \frac{1}{n} \sum_{i=1}^{n} h(\boldsymbol{x}, \boldsymbol{\omega}^i)$ be an epigraph variable

## Second Strategy: Decompose the Problem

$$\min_{\mathbf{x} \in \mathbb{R}^n, \theta} \quad \mathbf{c}^\top \mathbf{x} + \theta$$
$$\text{s.t.} \quad \mathbf{A}\mathbf{x} \leq \mathbf{b}.$$

## Second Strategy: Decompose the Problem

$$\min_{\mathbf{x} \in \mathbb{R}^n, \theta} \quad \mathbf{c}^\top \mathbf{x} + \theta$$
$$\text{s.t.} \quad \mathbf{A}\mathbf{x} \leq \mathbf{b}.$$

(Sketch) We iteratively

• Solve this "master" problem to find an optimal $\mathbf{x}$

## Second Strategy: Decompose the Problem

$$\min_{\boldsymbol{x} \in \mathbb{R}^n, \theta} \quad \boldsymbol{c}^\top \boldsymbol{x} + \theta$$
$$\text{s.t.} \quad \boldsymbol{A}\boldsymbol{x} \leq \boldsymbol{b}.$$

(Sketch) We iteratively

- Solve this "master" problem to find an optimal $\boldsymbol{x}$
- Evaluate $1/n \sum_{i=1}^{n} h(\boldsymbol{x}, \boldsymbol{\omega}^i)$ and add inequalities which model
  - $\theta \geq \frac{1}{n} \sum_{i=1}^{n} h(\boldsymbol{x}, \boldsymbol{\omega}^i)$
  - For $\boldsymbol{x}$ to be feasible, there is a feasible $\boldsymbol{y}(\boldsymbol{\omega}^i)$ in each scenario $\boldsymbol{\omega}^i$

until we converge.

## Second Strategy: Decompose the Problem

$$\min_{\boldsymbol{x} \in \mathbb{R}^n, \theta} \quad \boldsymbol{c}^\top \boldsymbol{x} + \theta$$

$$\text{s.t.} \quad \boldsymbol{A}\boldsymbol{x} \leq \boldsymbol{b}.$$

(Sketch) We iteratively

- Solve this "master" problem to find an optimal $\boldsymbol{x}$
- Evaluate $1/n \sum_{i=1}^n h(\boldsymbol{x}, \boldsymbol{\omega}^i)$ and add inequalities which model
  - $\theta \geq \frac{1}{n} \sum_{i=1}^n h(\boldsymbol{x}, \boldsymbol{\omega}^i)$
  - For $\boldsymbol{x}$ to be feasible, there is a feasible $\boldsymbol{y}(\boldsymbol{\omega}^i)$ in each scenario $\boldsymbol{\omega}^i$

until we converge. We never model $\boldsymbol{y}(\boldsymbol{\omega}^i)$, so we replaced one intractable problem with a sequence of (possibly many) tractable ones

## Second Strategy: Decompose the Problem

$$\min_{\mathbf{x} \in \mathbb{R}^n, \theta} \quad \mathbf{c}^\top \mathbf{x} + \theta$$

$$\text{s.t.} \quad \mathbf{A}\mathbf{x} \leq \mathbf{b}.$$

(Sketch) We iteratively

- Solve this "master" problem to find an optimal $\mathbf{x}$
- Evaluate $1/n \sum_{i=1}^n h(\mathbf{x}, \omega^i)$ and add inequalities which model
  - $\theta \geq \frac{1}{n} \sum_{i=1}^n h(\mathbf{x}, \omega^i)$
  - For $\mathbf{x}$ to be feasible, there is a feasible $\mathbf{y}(\omega^i)$ in each scenario $\omega^i$

until we converge. We never model $\mathbf{y}(\omega^i)$, so we replaced one intractable problem with a sequence of (possibly many) tractable ones

Remark: About to go through how this works in gory detail. However, I find the best way to understand this method is to code it for yourself.

## Benders Decomposition

Suppose we solve

$$\min_{\boldsymbol{x} \in \mathbb{R}^n, \theta} \quad \boldsymbol{c}^\top \boldsymbol{x} + \theta$$
$$\text{s.t.} \quad \boldsymbol{A}\boldsymbol{x} \leq \boldsymbol{b}.$$

and obtain some solution $\boldsymbol{x}$. Two cases:

- There is some scenario $\omega^i$ for which no $\boldsymbol{y}(\boldsymbol{\omega})$ can make the scenario feasible $\rightarrow$ we need to tell the master problem that this $\boldsymbol{x}$ is infeasible, via a *feasibility cut*

## Benders Decomposition

Suppose we solve

$$\min_{\boldsymbol{x} \in \mathbb{R}^n, \theta} \quad \boldsymbol{c}^\top \boldsymbol{x} + \theta$$
$$\text{s.t.} \quad \boldsymbol{A}\boldsymbol{x} \leq \boldsymbol{b}.$$

and obtain some solution $\boldsymbol{x}$. Two cases:

- There is some scenario $\omega^i$ for which no $\boldsymbol{y}(\boldsymbol{\omega})$ can make the scenario feasible $\rightarrow$ we need to tell the master problem that this $\boldsymbol{x}$ is infeasible, via a *feasibility cut*

- Every scenario $\omega^i$ is feasible $\rightarrow$ we need to tell the master problem how much $\boldsymbol{x}$ costs via an *optimality cut*

## Benders Decomposition: Feasibility Cut

Suppose we solve

$$\min_{\boldsymbol{x} \in \mathbb{R}^n, \theta} \quad \boldsymbol{c}^\top \boldsymbol{x} + \theta$$
$$\text{s.t.} \quad \boldsymbol{A}\boldsymbol{x} \leq \boldsymbol{b}.$$

and obtain some solution $\boldsymbol{x}$ such that in scenario $i$ no $\boldsymbol{y}(\boldsymbol{\omega})$ can make the scenario feasible.

## Benders Decomposition: Feasibility Cut

Suppose we solve

$$\min_{\boldsymbol{x}\in\mathbb{R}^n,\theta} \quad \boldsymbol{c}^\top\boldsymbol{x} + \theta$$
$$\text{s.t.} \quad \boldsymbol{A}\boldsymbol{x} \leq \boldsymbol{b}.$$

and obtain some solution $\boldsymbol{x}$ such that in scenario $i$ no $\boldsymbol{y}(\boldsymbol{\omega})$ can make the scenario feasible. Then, the dual problem in this scenario is unbounded (why?), so there is some $\boldsymbol{\mu}(\boldsymbol{\omega}^i)$ such that

$$(\boldsymbol{d}(\omega) - \boldsymbol{D}(\omega)\boldsymbol{x})^\top\boldsymbol{\mu}(\omega) > 0, \ \boldsymbol{F}(\omega)^\top\boldsymbol{\mu}(\omega) = \boldsymbol{q}(\omega), \boldsymbol{\mu}(\omega) \leq \boldsymbol{0}.$$

## Benders Decomposition: Feasibility Cut

Suppose we solve

$$\min_{\mathbf{x}\in\mathbb{R}^n,\theta} \quad \mathbf{c}^\top\mathbf{x} + \theta$$
$$\text{s.t.} \quad \mathbf{A}\mathbf{x} \leq \mathbf{b}.$$

and obtain some solution $\mathbf{x}$ such that in scenario $i$ no $\mathbf{y}(\omega)$ can make the scenario feasible. Then, the dual problem in this scenario is unbounded (why?), so there is some $\boldsymbol{\mu}(\omega^i)$ such that

$$(\mathbf{d}(\omega) - \mathbf{D}(\omega)\mathbf{x})^\top\boldsymbol{\mu}(\omega) > 0, \ \mathbf{F}(\omega)^\top\boldsymbol{\mu}(\omega) = \mathbf{q}(\omega), \boldsymbol{\mu}(\omega) \leq \mathbf{0}.$$

Therefore, we fix $\boldsymbol{\mu}(\omega^i)$ and impose the feasibility cut

$$(\mathbf{d}(\omega^i) - \mathbf{D}(\omega^i)\mathbf{x})^\top\boldsymbol{\mu}(\omega^i) \leq 0,$$

in the master problem, where everything but $\mathbf{x}$ is data

## In This Case, The Master Problem Now Looks Like

$$\min_{\boldsymbol{x} \in \mathbb{R}^n, \theta} \quad \boldsymbol{c}^\top \boldsymbol{x} + \theta$$
$$\text{s.t.} \quad \boldsymbol{A}\boldsymbol{x} \leq \boldsymbol{b},$$
$$(\boldsymbol{d}(\boldsymbol{\omega}^i) - \boldsymbol{D}(\boldsymbol{\omega}^i)\boldsymbol{x})^\top \boldsymbol{\mu}(\boldsymbol{\omega}^i) \leq 0.$$

## Case Two: Each Scenario Was Feasible

$\theta$ has usually underestimated $1/n \sum_{i=1}^{n} h(\boldsymbol{x}, \boldsymbol{\omega}^i)$

## Case Two: Each Scenario Was Feasible

$\theta$ has usually underestimated $1/n \sum_{i=1}^{n} h(\mathbf{x}, \boldsymbol{\omega}^i)$

Need cut involving $\theta$, which tells master problem what $\mathbf{x}$ costs

## Case Two: Each Scenario Was Feasible

$\theta$ has usually underestimated $1/n \sum_{i=1}^{n} h(\mathbf{x}, \boldsymbol{\omega}^i)$

Need cut involving $\theta$, which tells master problem what $\mathbf{x}$ costs

By strong duality

$$\frac{1}{n} \sum_{i=1}^{n} h(\mathbf{x}, \boldsymbol{\omega}^i) = 1/n \sum_{i=1}^{n} (\mathbf{d}(\boldsymbol{\omega}^i) - \mathbf{D}(\boldsymbol{\omega}^i)\mathbf{x})^{\top} \boldsymbol{\mu}(\boldsymbol{\omega}^i),$$

where $\boldsymbol{\mu}(\boldsymbol{\omega}^i)$, dual-optimal in scenario $i$, is data

## Case Two: Each Scenario Was Feasible

$\theta$ has usually underestimated $1/n \sum_{i=1}^{n} h(\boldsymbol{x}, \boldsymbol{\omega}^i)$

Need cut involving $\theta$, which tells master problem what $\boldsymbol{x}$ costs

By strong duality

$$\frac{1}{n} \sum_{i=1}^{n} h(\boldsymbol{x}, \boldsymbol{\omega}^i) = 1/n \sum_{i=1}^{n} (\boldsymbol{d}(\boldsymbol{\omega}^i) - \boldsymbol{D}(\boldsymbol{\omega}^i)\boldsymbol{x})^{\top} \boldsymbol{\mu}(\boldsymbol{\omega}^i),$$

where $\boldsymbol{\mu}(\boldsymbol{\omega}^i)$, dual-optimal in scenario $i$, is data

By weak duality, for any $\bar{\boldsymbol{x}}$

$$\frac{1}{n} \sum_{i=1}^{n} h(\bar{\boldsymbol{x}}, \boldsymbol{\omega}^i) \geq \frac{1}{n} \sum_{i=1}^{n} (\boldsymbol{d}(\boldsymbol{\omega}^i) - \boldsymbol{D}(\boldsymbol{\omega}^i)\bar{\boldsymbol{x}})^{\top} \boldsymbol{\mu}(\boldsymbol{\omega}^i),$$

where everything but $\bar{\boldsymbol{x}}$ is data

## Case Two: Each Scenario Was Feasible

$\theta$ has usually underestimated $1/n \sum_{i=1}^{n} h(\boldsymbol{x}, \boldsymbol{\omega}^i)$

Need cut involving $\theta$, which tells master problem what $\boldsymbol{x}$ costs

By strong duality

$$\frac{1}{n} \sum_{i=1}^{n} h(\boldsymbol{x}, \boldsymbol{\omega}^i) = 1/n \sum_{i=1}^{n} (\boldsymbol{d}(\boldsymbol{\omega}^i) - \boldsymbol{D}(\boldsymbol{\omega}^i)\boldsymbol{x})^{\top} \boldsymbol{\mu}(\boldsymbol{\omega}^i),$$

where $\boldsymbol{\mu}(\boldsymbol{\omega}^i)$, dual-optimal in scenario $i$, is data

By weak duality, for any $\bar{\boldsymbol{x}}$

$$\frac{1}{n} \sum_{i=1}^{n} h(\bar{\boldsymbol{x}}, \boldsymbol{\omega}^i) \geq \frac{1}{n} \sum_{i=1}^{n} (\boldsymbol{d}(\boldsymbol{\omega}^i) - \boldsymbol{D}(\boldsymbol{\omega}^i)\bar{\boldsymbol{x}})^{\top} \boldsymbol{\mu}(\boldsymbol{\omega}^i),$$

where everything but $\bar{\boldsymbol{x}}$ is data

Therefore, we add cut

$$\theta \geq \frac{1}{n} \sum_{i=1}^{n} (\boldsymbol{d}(\boldsymbol{\omega}^i) - \boldsymbol{D}(\boldsymbol{\omega}^i)\bar{\boldsymbol{x}})^{\top} \boldsymbol{\mu}(\boldsymbol{\omega}^i)$$

28

## The Master Problem Might Now Look Like

$$\min_{\mathbf{x}\in\mathbb{R}^n,\theta} \quad \mathbf{c}^\top\mathbf{x} + \theta$$

$$\text{s.t.} \quad \mathbf{A}\mathbf{x} \leq \mathbf{b},$$

$$\theta \geq \frac{1}{n}\sum_{i=1}^{n}(\mathbf{d}(\omega^i) - \mathbf{D}(\omega^i)\bar{\mathbf{x}})^\top\boldsymbol{\mu}(\omega^i),$$

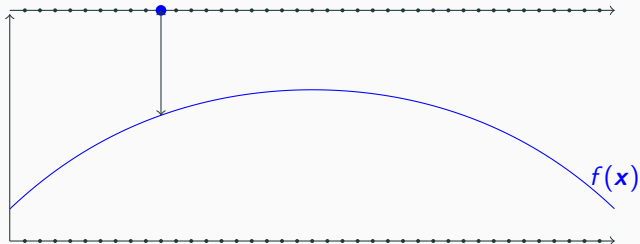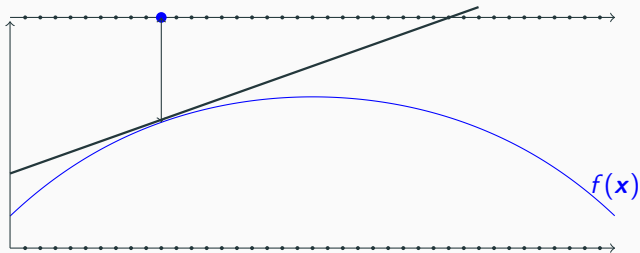$$(\mathbf{d}(\omega^i) - \mathbf{D}(\omega^i)\mathbf{x})^\top\boldsymbol{\mu}(\omega^i) \leq 0.$$

$f(x)$

$f(\boldsymbol{x})$

$f(\boldsymbol{x})$

$f(\mathbf{x})$

$f(\boldsymbol{x})$

$f(\mathbf{x})$

**Sample Average Approximation: Code
You will write this yourself in the first
assignment :-)**

# Can we do Better? Ridge Regression and Sample-Average Approximation

# Can we do Better Than the Sample-Average Approximation?

## Returning to Linear Regression

Statisticians don't solve problems like

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^p} \quad \frac{1}{n}\|\boldsymbol{X}\boldsymbol{\beta} - \boldsymbol{y}\|_2^2$$

to pick $\boldsymbol{\beta}$, despite SAA's properties. Why not?

## Returning to Linear Regression

Statisticians don't solve problems like

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^p} \quad \frac{1}{n} \|\boldsymbol{X}\boldsymbol{\beta} - \boldsymbol{y}\|_2^2$$

to pick $\boldsymbol{\beta}$, despite SAA's properties. Why not?

Because $n$ is finite; we want $\boldsymbol{\beta}$ to perform as well as possible on an unseen observation $(\boldsymbol{x}_i, y_i)$, not just minimize training error.

**Returning to Linear Regression**

Statisticians don't solve problems like

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^p} \quad \frac{1}{n} \|\boldsymbol{X}\boldsymbol{\beta} - \boldsymbol{y}\|_2^2$$

to pick $\boldsymbol{\beta}$, despite SAA's properties. Why not?

Because $n$ is finite; we want $\boldsymbol{\beta}$ to perform as well as possible on an unseen observation $(\boldsymbol{x}_i, y_i)$, not just minimize training error. They solve

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^p} \quad \frac{1}{n} \|\boldsymbol{X}\boldsymbol{\beta} - \boldsymbol{y}\|_2^2 + R(\boldsymbol{\beta}),$$

where $R(\cdot)$ is a regularization term, e.g., $\frac{1}{2\gamma}\|\boldsymbol{\beta}\|_2^2 + \lambda\|\boldsymbol{\beta}\|_1$ for appropriately chosen $\lambda, \gamma$ (elastic net method, Zou and Hastie 2005).

## Returning to Linear Regression

Statisticians don't solve problems like

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^p} \quad \frac{1}{n} \|\boldsymbol{X}\boldsymbol{\beta} - \boldsymbol{y}\|_2^2$$

to pick $\boldsymbol{\beta}$, despite SAA's properties. Why not?

Because $n$ is finite; we want $\boldsymbol{\beta}$ to perform as well as possible on an unseen observation $(\boldsymbol{x}_i, y_i)$, not just minimize training error. They solve

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^p} \quad \frac{1}{n} \|\boldsymbol{X}\boldsymbol{\beta} - \boldsymbol{y}\|_2^2 + R(\boldsymbol{\beta}),$$

where $R(\cdot)$ is a regularization term, e.g., $\frac{1}{2\gamma}\|\boldsymbol{\beta}\|_2^2 + \lambda\|\boldsymbol{\beta}\|_1$ for appropriately chosen $\lambda, \gamma$ (elastic net method, Zou and Hastie 2005).

This usually performs better out-of-sample.

## What's That Gotta Do With The Price of Fish?

- In the 2000s, sample-average approximation was a very popular method for optimizing under uncertainty

## What's That Gotta Do With The Price of Fish?

- In the 2000s, sample-average approximation was a very popular method for optimizing under uncertainty

- In the early 2010s, the community became more aware of the danger of overfitting. Since then, variants of SAA that account for overfitting with better finite-sample guarantees have become popular

## What's That Gotta Do With The Price of Fish?

- In the 2000s, sample-average approximation was a very popular method for optimizing under uncertainty
- In the early 2010s, the community became more aware of the danger of overfitting. Since then, variants of SAA that account for overfitting with better finite-sample guarantees have become popular
- We still teach SAA, because you need to understand SAA first

## What's That Gotta Do With The Price of Fish?

- In the 2000s, sample-average approximation was a very popular method for optimizing under uncertainty
- In the early 2010s, the community became more aware of the danger of overfitting. Since then, variants of SAA that account for overfitting with better finite-sample guarantees have become popular
- We still teach SAA, because you need to understand SAA first
- Variants intimately related to distributional robustness, so we'll come back to them later

## What's That Gotta Do With The Price of Fish?

- In the 2000s, sample-average approximation was a very popular method for optimizing under uncertainty
- In the early 2010s, the community became more aware of the danger of overfitting. Since then, variants of SAA that account for overfitting with better finite-sample guarantees have become popular
- We still teach SAA, because you need to understand SAA first
- Variants intimately related to distributional robustness, so we'll come back to them later
- Google "Robust SAA" by Bertsimas et al. (Math. Prog. 2017)

**Extension: Accelerating Benders Decomposition for Facility Location**

See slides by Fischetti (2017)

# Activities for if we Finish Early

## Either Prove or Provide a Counterexample for the Following Statements

- The intersection of convex sets is convex.
- The union of convex sets is convex.
- All polyhedral sets are convex.

## Some (Classically) Useful Terms

- Value of Stochastic Solution.
- Value of Perfect Information.

## (More) Activities for if we Finish Early

1. HW1 Q0.

## (More) Activities for if we Finish Early

1. HW1 Q0.
2. Class discussion: Summarize the Pros and Cons of the Sample Average Approximation Method, based on what we have learned so far.

## (More) Activities for if we Finish Early

1. HW1 Q0.
2. Class discussion: Summarize the Pros and Cons of the Sample Average Approximation Method, based on what we have learned so far.
3. Shapiro and Philpott Introduction to Stochastic Programming Tutorial.

## (More) Activities for if we Finish Early

1. HW1 Q0.
2. Class discussion: Summarize the Pros and Cons of the Sample Average Approximation Method, based on what we have learned so far.
3. Shapiro and Philpott Introduction to Stochastic Programming Tutorial.
4. Open office hours.

**Suggested Readings**

## Suggested Readings to Accompany Today's Lecture

A friendly reminder:

> *"To get as much out of this class as possible, we suggest that you spend at least as much time on reading the papers and textbooks referenced in the lectures/reviewing the lectures as you spend in class."* — *The syllabus*

## Suggested Readings to Accompany Today's Lecture

A friendly reminder:

> *"To get as much out of this class as possible, we suggest that you spend at least as much time on reading the papers and textbooks referenced in the lectures/reviewing the lectures as you spend in class."* — The syllabus

Recommended reading:

- Shapiro, Dentcheva, Ruszczynski *Lectures on Stochastic Programming: Modeling and Theory* (2013), Chapters 1.1 and 2.

Optional further reading:

- Recht *Lecture* 1. In CS294 The Mathematics of Data Science lecture notes, UC Berkeley (2013).
- Kim, Pasupathy, Henderson *A Guide to Sample-Average Approximation*. In: Handbook of simulation optimization (2015).
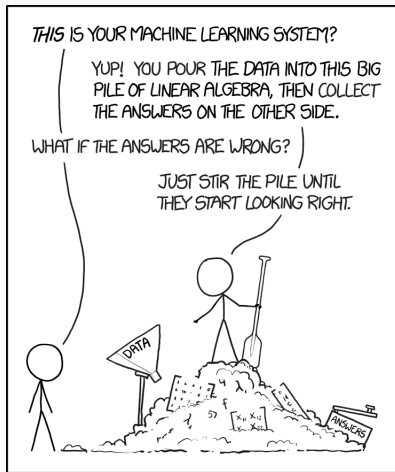
**Figure 3:** There's *always* a relevant XKCD

**Thank you, and see you next week!**