

6.7900: Machine Learning Lecture 6

Lecture start: Tues/Thurs 2:35pm

Who's speaking today? Prof. Tamara Broderick

Course website: gradml.mit.edu

Questions? Ask here or on piazza.com/mit/fall2024/67900/

Materials: Slides, video, etc linked from gradml.mit.edu after the lecture (but there is no livestream)

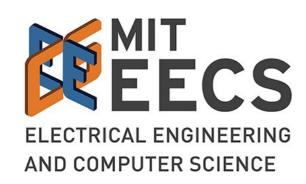
Last Time

Challenges with MLE/ERM I. Visualizations for linear regression

- II. Bayesian linear regression
 - A. Bayes with Gaussians
 - B. Posterior

Today

- II. Uncertainty
- III. Ridge regression
- IV. More flexible/complex features



6.7900: Machine Learning Lecture 6

Lecture start: Tues/Thurs 2:35pm

Who's speaking today? Prof. Tamara Broderick

Course website: gradml.mit.edu

Questions? Ask here or on piazza.com/mit/fall2024/67900/

Materials: Slides, video, etc linked from gradml.mit.edu after the lecture (but there is no livestream)

Last Time

- Challenges with MLE/ERM I. Visualizations for linear regression
- II. Bayesian linear regression
 - A. Bayes with Gaussians
 - B. Posterior

Today

- II. Uncertainty
- III. Ridge regression
- IV. More flexible/complex features

See board and demos

 Rather than return an error in the face of collinearity, the Bayesian approach expresses uncertainty over possibilities

- Rather than return an error in the face of collinearity, the Bayesian approach expresses uncertainty over possibilities
- It's common to see publications discuss:

- Rather than return an error in the face of collinearity, the Bayesian approach expresses uncertainty over possibilities
- It's common to see publications discuss: [Aldous, "Is it practical and useful to distinguish
 - Aleatoric uncertainty: "intrinsic randomness" aleatoric and epistemic uncertainty?"]

- Rather than return an error in the face of collinearity, the Bayesian approach expresses uncertainty over possibilities
- It's common to see publications discuss: [Aldous, "Is it practical and useful to distinguish
 - Aleatoric uncertainty: "intrinsic randomness" aleatoric and epistemic uncertainty?"]
 - Epistemic uncertainty: "lack of knowledge" about something that we could know for sure in principle

- Rather than return an error in the face of collinearity, the Bayesian approach expresses uncertainty over possibilities
- It's common to see publications discuss: [Aldous, "Is it practical and useful to distinguish
 - Aleatoric uncertainty: "intrinsic randomness" aleatoric and epistemic uncertainty?"]
 - Epistemic uncertainty: "lack of knowledge" about something that we could know for sure in principle
 - Example: I plan to shuffle a deck of cards and then ask whether the top card is an ace.

- Rather than return an error in the face of collinearity, the Bayesian approach expresses uncertainty over possibilities
- It's common to see publications discuss: [Aldous, "Is it practical and useful to distinguish
 - Aleatoric uncertainty: "intrinsic randomness" aleatoric and epistemic uncertainty?"]
 - Epistemic uncertainty: "lack of knowledge" about something that we could know for sure in principle
 - Example: I plan to shuffle a deck of cards and then ask whether the top card is an ace.
 - Before I shuffle, uncertainty is aleatoric. After shuffling but before looking, uncertainty is epistemic.

- Rather than return an error in the face of collinearity, the Bayesian approach expresses uncertainty over possibilities
- It's common to see publications discuss: [Aldous, "Is it practical and useful to distinguish
 - Aleatoric uncertainty: "intrinsic randomness" aleatoric and epistemic uncertainty?"]
 - Epistemic uncertainty: "lack of knowledge" about something that we could know for sure in principle
 - Example: I plan to shuffle a deck of cards and then ask whether the top card is an ace.
 - Before I shuffle, uncertainty is aleatoric. After shuffling but before looking, uncertainty is epistemic.

 Predictive vs posterior

- Rather than return an error in the face of collinearity, the Bayesian approach expresses uncertainty over possibilities
- It's common to see publications discuss:
 - Aleatoric uncertainty: "intrinsic randomness" aleatoric and epistemic uncertainty?"]

[Aldous, "Is it practical

- Epistemic uncertainty: "lack of knowledge" about something that we could know for sure in principle
- Example: I plan to shuffle a deck of cards and then ask whether the top card is an ace.
 - Before I shuffle, uncertainty is aleatoric. After shuffling but before looking, uncertainty is epistemic. Predictive vs posterior

- Rather than return an error in the face of collinearity, the Bayesian approach expresses uncertainty over possibilities
- It's common to see publications discuss: [Aldous, "Is it practical and useful to distinguish
 - Aleatoric uncertainty: "intrinsic randomness" aleatoric and epistemic uncertainty?"]
 - Epistemic uncertainty: "lack of knowledge" about something that we could know for sure in principle
 - Example: I plan to shuffle a deck of cards and then ask whether the top card is an ace.
 - Before I shuffle, uncertainty is aleatoric. After shuffling but before looking, uncertainty is epistemic.
- The Bayesian posterior uncertainty is conditional on our model. Goes to zero if enough data with enough x coverage.

- Rather than return an error in the face of collinearity, the Bayesian approach expresses uncertainty over possibilities
- It's common to see publications discuss: [Aldous, "Is it practical and useful to distinguish
 - Aleatoric uncertainty: "intrinsic randomness" aleatoric and epistemic uncertainty?"]
 - Epistemic uncertainty: "lack of knowledge" about something that we could know for sure in principle
 - Example: I plan to shuffle a deck of cards and then ask whether the top card is an ace.
 - Before I shuffle, uncertainty is aleatoric. After shuffling but before looking, uncertainty is epistemic.
- The Bayesian posterior uncertainty is conditional on our model. Goes to zero if enough data with enough x coverage.
 - "Uncertainty" as a term of art vs colloquial meaning

- Rather than return an error in the face of collinearity, the Bayesian approach expresses uncertainty over possibilities
- It's common to see publications discuss: [Aldous, "Is it practical and useful to distinguish
 - Aleatoric uncertainty: "intrinsic randomness" aleatoric and epistemic uncertainty?"]
 - Epistemic uncertainty: "lack of knowledge" about something that we could know for sure in principle
 - Example: I plan to shuffle a deck of cards and then ask whether the top card is an ace.
 - Before I shuffle, uncertainty is aleatoric. After shuffling but before looking, uncertainty is epistemic.
- The Bayesian posterior uncertainty is conditional on our model. Goes to zero if enough data with enough x coverage.
 - "Uncertainty" as a term of art vs colloquial meaning
 - Not just a Bayesian issue: (frequentist) standard errors go to 0 as data grows. So in real-life ("all models are wrong"), common p-values get arbitrarily small from enough data

 What if I just want a single regression line and not a whole distribution?

• What if I just want a single regression line and not a whole distribution? Some options: posterior mean or MAP.

 What if I just want a single regression line and not a whole distribution? Some options: posterior mean or MAP. Lecture 3:

 What if I just want a single regression line and not a whole distribution? Some options: posterior mean or MAP. Lecture 3:

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} p(\theta|\mathcal{D})$$

 What if I just want a single regression line and not a whole distribution? Some options: posterior mean or MAP. Lecture 3:

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} p(\theta|\mathcal{D}) = \arg \min_{\theta} - \log p(\theta|\mathcal{D})$$

 What if I just want a single regression line and not a whole distribution? Some options: posterior mean or MAP. Lecture 3:

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} p(\theta|\mathcal{D}) = \arg \min_{\theta} - \log p(\theta|\mathcal{D})$$

• consider the linear regression model we used, now with $\mu_0=0_D$ & $\Sigma_0=\sigma_0^2I_{D\times D}$

 What if I just want a single regression line and not a whole distribution? Some options: posterior mean or MAP. Lecture 3:

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} p(\theta|\mathcal{D}) = \arg \min_{\theta} - \log p(\theta|\mathcal{D})$$

• consider the linear regression model we used, now with $\mu_0=0_D$ & $\Sigma_0=\sigma_0^2I_{D\times D}$

$$\begin{aligned} &\text{now with } \mu_0 = 0_D \quad \& \quad \Sigma_0 = \sigma_0^2 I_{D \times D} \\ &= \arg\min_{\theta} \left\{ (X\theta - Y)^\top (X\theta - Y) + \frac{\sigma^2}{\sigma_0^2} \theta^\top \theta \right\} \end{aligned}$$

 What if I just want a single regression line and not a whole distribution? Some options: posterior mean or MAP. Lecture 3:

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} p(\theta|\mathcal{D}) = \arg \min_{\theta} - \log p(\theta|\mathcal{D})$$

• consider the linear regression model we used,

$$\begin{aligned} &\text{now with } \mu_0 = 0_D \quad \& \quad \Sigma_0 = \sigma_0^2 I_{D \times D} \\ &= \arg\min_{\theta} \left\{ (X\theta - Y)^\top (X\theta - Y) + \frac{\sigma^2}{\sigma_0^2} \theta^\top \theta \right\} \end{aligned}$$

the usual MLE/ERM objective

 What if I just want a single regression line and not a whole distribution? Some options: posterior mean or MAP. Lecture 3:

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} p(\theta|\mathcal{D}) = \arg \min_{\theta} - \log p(\theta|\mathcal{D})$$

• consider the linear regression model we used,

$$\begin{aligned} & \text{now with } \ \mu_0 = 0_D \quad \& \quad \Sigma_0 = \sigma_0^2 I_{D \times D} \\ &= \arg\min_{\theta} \left\{ (X\theta - Y)^\top (X\theta - Y) + \frac{\sigma^2}{\sigma_0^2} \theta^\top \theta \right\} \end{aligned}$$

 What if I just want a single regression line and not a whole distribution? Some options: posterior mean or MAP. Lecture 3:

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} p(\theta|\mathcal{D}) = \arg \min_{\theta} - \log p(\theta|\mathcal{D})$$

• consider the linear regression model we used, now with $\mu_0=0_D$ & $\Sigma_0=\sigma_0^2I_{D\times D}$

$$\begin{array}{ll} \text{ridge} \\ \text{regression} \end{array} = \arg\min_{\theta} \left\{ (X\theta - Y)^{\top} (X\theta - Y) + \frac{\sigma^2}{\sigma_0^2} \theta^{\top} \theta \right\}$$

the usual MLE/ERM (2/ridge) penalty/ objective

regularizer

 What if I just want a single regression line and not a whole distribution? Some options: posterior mean or MAP. Lecture 3:

$$\hat{\theta}_{\text{MAP}} = \arg\max_{\theta} p(\theta|\mathcal{D}) = \arg\min_{\theta} - \log p(\theta|\mathcal{D})$$

• consider the linear regression model we used, now with $\,\mu_0=0_D\,$ & $\,\Sigma_0=\sigma_0^2I_{D\times D}\,$

$$= \arg\min_{\theta} \left\{ (X\theta - Y)^{\top} (X\theta - Y) + \frac{\sigma^2}{\sigma_0^2} \theta^{\top} \theta \right\}$$

["ridge analysis" RW Hoerl 2020, Ridge Regression: A Historical Context; can think of "ridge" here due to collinearity in likelihood]

the usual MLE/ERM objective

(\ell2/ridge) penalty/ regularizer

 What if I just want a single regression line and not a whole distribution? Some options: posterior mean or MAP. Lecture 3:

$$\hat{\theta}_{\text{MAP}} = \arg\max_{\theta} p(\theta|\mathcal{D}) = \arg\min_{\theta} - \log p(\theta|\mathcal{D})$$

• consider the linear regression model we used,

$$\text{now with } \mu_0 = 0_D \quad \& \quad \Sigma_0 = \sigma_0^2 I_{D \times D}$$

$$\text{ridge regression} = \arg\min_{\theta} \left\{ (X\theta - Y)^\top (X\theta - Y) + \frac{\sigma^2}{\sigma_0^2} \theta^\top \theta \right\}$$
 ["ridge analysis" RW Hoerl 2020,

Ridge Regression: A Historical Context; can think of "ridge" here

due to collinearity in likelihood]

the usual MLE/ERM (l2/ridge) penalty/ regularizer

 Since the Gaussian posterior is symmetric, the MAP is equal to the mean here. So we already computed it in Lecture 5.

 What if I just want a single regression line and not a whole distribution? Some options: posterior mean or MAP. Lecture 3:

$$\hat{\theta}_{\text{MAP}} = \arg\max_{\theta} p(\theta|\mathcal{D}) = \arg\min_{\theta} - \log p(\theta|\mathcal{D})$$

• consider the linear regression model we used, now with $\mu_0=0_D$ & $\Sigma_0=\sigma_0^2I_{D\times D}$

$$\begin{array}{ll} \textbf{ridge} \\ \textbf{regression} \\ \textbf{"ridge analysis" BW Hoerl 2020.} \end{array} = \arg\min_{\theta} \left\{ (X\theta - Y)^{\top} (X\theta - Y) + \frac{\sigma^2}{\sigma_0^2} \theta^{\top} \theta \right\}$$

["ridge analysis" RW Hoerl 2020, Ridge Regression: A Historical Context; can think of "ridge" here due to collinearity in likelihood]

- Since the Gaussian posterior is symmetric, the MAP is equal to the mean here. So we already computed it in Lecture 5.
 - Letting $\lambda = \sigma^2/\sigma_0^2$,

 What if I just want a single regression line and not a whole distribution? Some options: posterior mean or MAP. Lecture 3:

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} p(\theta|\mathcal{D}) = \arg \min_{\theta} - \log p(\theta|\mathcal{D})$$

• consider the linear regression model we used,

$$\text{now with } \mu_0 = 0_D \quad \& \quad \Sigma_0 = \sigma_0^2 I_{D \times D}$$

$$\underset{\text{ridge analysis" BW Hoerl 2020}}{\text{ridge analysis" BW Hoerl 2020}} = \arg\min_{\theta} \left\{ (X\theta - Y)^\top (X\theta - Y) + \frac{\sigma^2}{\sigma_0^2} \theta^\top \theta \right\}$$

["ridge analysis" RW Hoerl 2020, Ridge Regression: A Historical

Context; can think of "ridge" here due to collinearity in likelihood]

- Since the Gaussian posterior is symmetric, the MAP is equal to the mean here. So we already computed it in Lecture 5.
 - Letting $\lambda = \sigma^2/\sigma_0^2$, $\hat{\theta}_{\mathrm{MAP}} = (\lambda I_{D\times D} + X^{\top}X)^{-1}(X^{\top}Y)$

 What if I just want a single regression line and not a whole distribution? Some options: posterior mean or MAP. Lecture 3:

$$\hat{\theta}_{\text{MAP}} = \arg\max_{\theta} p(\theta|\mathcal{D}) = \arg\min_{\theta} - \log p(\theta|\mathcal{D})$$

• consider the linear regression model we used, now with $\mu_0=0_D$ & $\Sigma_0=\sigma_0^2I_{D\times D}$

$$\begin{array}{ll} \text{ridge} \\ \text{regression} \\ \text{["ridge analysis" RW Hoerl 2020,} \end{array} = \arg\min_{\theta} \left\{ (X\theta - Y)^{\top} (X\theta - Y) + \frac{\sigma^2}{\sigma_0^2} \theta^{\top} \theta \right\}$$

Ridge Regression: A Historical Context; can think of "ridge" here due to collinearity in likelihood]

the usual MLE/ERM (l2/ridge) penalty/objective regularizer

- Since the Gaussian posterior is symmetric, the MAP is equal to the mean here. So we already computed it in Lecture 5.
 - Letting $\lambda = \sigma^2/\sigma_0^2$, $\hat{\theta}_{\text{MAP}} = (\lambda I_{D\times D} + X^{\top}X)^{-1}(X^{\top}Y)$

What happens as λ gets very small?

• What if I just want a single regression line and not a whole distribution? Some options: posterior mean or MAP. Lecture 3:

$$\hat{\theta}_{\text{MAP}} = \arg\max_{\theta} p(\theta|\mathcal{D}) = \arg\min_{\theta} - \log p(\theta|\mathcal{D})$$

• consider the linear regression model we used, now with $\mu_0=0_D$ & $\Sigma_0=\sigma_0^2I_{D\times D}$

$$\begin{array}{ll} \text{ridge} \\ \text{regression} \\ \text{["ridge analysis" RW Hoerl 2020,} \end{array} = \arg\min_{\theta} \left\{ (X\theta - Y)^{\top} (X\theta - Y) + \frac{\sigma^2}{\sigma_0^2} \theta^{\top} \theta \right\}$$

Ridge Regression: A Historical Context; can think of "ridge" here due to collinearity in likelihood]

the usual MLE/ERM (l2/ridge) penalty/objective regularizer

- Since the Gaussian posterior is symmetric, the MAP is equal to the mean here. So we already computed it in Lecture 5.
 - Letting $\lambda = \sigma^2/\sigma_0^2$, $\hat{\theta}_{\text{MAP}} = (\lambda I_{D\times D} + X^{\top}X)^{-1}(X^{\top}Y)$

What happens as λ gets very small? Very large?

• What if I just want a single regression line and not a whole distribution? Some options: posterior mean or MAP. Lecture 3:

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} p(\theta|\mathcal{D}) = \arg \min_{\theta} - \log p(\theta|\mathcal{D})$$

• consider the linear regression model we used, now with $\mu_0=0_D$ & $\Sigma_0=\sigma_0^2I_{D\times D}$

$$\begin{array}{ll} \textbf{ridge} \\ \textbf{regression} \\ \textbf{regression} \\ \textbf{fridge analysis" BW Hoerl 2020} \end{array} = \arg\min_{\theta} \left\{ (X\theta - Y)^{\top} (X\theta - Y) + \frac{\sigma^2}{\sigma_0^2} \theta^{\top} \theta \right\}$$

["ridge analysis" RW Hoerl 2020, Ridge Regression: A Historical Context; can think of "ridge" here due to collinearity in likelihood]

- Since the Gaussian posterior is symmetric, the MAP is equal to the mean here. So we already computed it in Lecture 5.
 - Letting $\lambda = \sigma^2/\sigma_0^2$, $\hat{\theta}_{\text{MAP}} = (\lambda I_{D\times D} + X^{\top}X)^{-1}(X^{\top}Y)$
- Avoids the inversion issue of the MLE by choosing a particular option among options "equally good" in MLE

 What if I just want a single regression line and not a whole distribution? Some options: posterior mean or MAP. Lecture 3:

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} p(\theta|\mathcal{D}) = \arg \min_{\theta} - \log p(\theta|\mathcal{D})$$

• consider the linear regression model we used, now with $\mu_0=0_D$ & $\Sigma_0=\sigma_0^2I_{D\times D}$

$$\begin{array}{ll} \textbf{ridge} \\ \textbf{regression} \\ \textbf{regression} \\ \textbf{fridge analysis" BW Hoerl 2020} \end{array} = \arg\min_{\theta} \left\{ (X\theta - Y)^{\top} (X\theta - Y) + \frac{\sigma^2}{\sigma_0^2} \theta^{\top} \theta \right\}$$

["ridge analysis" RW Hoerl 2020, Ridge Regression: A Historical Context; can think of "ridge" here due to collinearity in likelihood]

- Since the Gaussian posterior is symmetric, the MAP is equal to the mean here. So we already computed it in Lecture 5.
 - Letting $\lambda = \sigma^2/\sigma_0^2$, $\hat{\theta}_{\text{MAP}} = (\lambda I_{D\times D} + X^{\top}X)^{-1}(X^{\top}Y)$
- Avoids the inversion issue of the MLE by choosing a particular option among options "equally good" in MLE check: always invertible

 What if I just want a single regression line and not a whole distribution? Some options: posterior mean or MAP. Lecture 3:

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} p(\theta|\mathcal{D}) = \arg \min_{\theta} - \log p(\theta|\mathcal{D})$$

• consider the linear regression model we used, now with $\mu_0=0_D$ & $\Sigma_0=\sigma_0^2I_{D\times D}$

$$\begin{array}{ll} \text{ridge} \\ \text{regression} \\ \text{["ridge analysis" RW Hoerl 2020,} \end{array} = \arg\min_{\theta} \left\{ (X\theta - Y)^{\top} (X\theta - Y) + \frac{\sigma^2}{\sigma_0^2} \theta^{\top} \theta \right\}$$

Ridge Regression: A Historical Context; can think of "ridge" here due to collinearity in likelihood]

- Since the Gaussian posterior is symmetric, the MAP is equal to the mean here. So we already computed it in Lecture 5.
 - Letting $\lambda = \sigma^2/\sigma_0^2$, $\hat{\theta}_{\text{MAP}} = (\lambda I_{D\times D} + X^{\top}X)^{-1}(X^{\top}Y)$
- Avoids the inversion issue of the MLE by choosing a particular option among options "equally good" in MLE
 - No uncertainty (e.g. when regression is underdetermined)

 What if I just want a single regression line and not a whole distribution? Some options: posterior mean or MAP. Lecture 3:

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} p(\theta|\mathcal{D}) = \arg \min_{\theta} - \log p(\theta|\mathcal{D})$$

• consider the linear regression model we used, now with $\mu_0=0_D$ & $\Sigma_0=\sigma_0^2I_{D\times D}$

$$\begin{array}{ll} \text{ridge} \\ \text{regression} \\ \text{["ridge analysis" RW Hoerl 2020,} \end{array} = \arg\min_{\theta} \left\{ (X\theta - Y)^{\top} (X\theta - Y) + \frac{\sigma^2}{\sigma_0^2} \theta^{\top} \theta \right\}$$

Ridge Regression: A Historical Context; can think of "ridge" here due to collinearity in likelihood]

the usual MLE/ERM (@/ridge) penalty/ objective

regularizer

- Since the Gaussian posterior is symmetric, the MAP is equal to the mean here. So we already computed it in Lecture 5.
 - Letting $\lambda = \sigma^2/\sigma_0^2$, $\hat{\theta}_{\text{MAP}} = (\lambda I_{D\times D} + X^{\top}X)^{-1}(X^{\top}Y)$
- Avoids the inversion issue of the MLE by choosing a particular option among options "equally good" in MLE
 - No uncertainty (e.g. when regression is underdetermined)
 - May offer better generalization than vanilla MLE

A note on features

A note on features

• Linear models can be very flexible given non-trivial features

- Linear models can be very flexible given non-trivial features
- We've been considering

- Linear models can be very flexible given non-trivial features
- We've been considering $h(x) = \theta_1 x_1 + \dots + \theta_D x_D = \theta^\top x$

- Linear models can be very flexible given non-trivial features
- We've been considering $h(x) = \theta_1 x_1 + \dots + \theta_D x_D = \theta^\top x$
- But we could take $h(x) = \hat{\theta}_1 \phi_1(x) + \dots + \theta_D \phi_D(x) = \theta^\top \phi(x)$

- Linear models can be very flexible given non-trivial features
- We've been considering $h(x) = \theta_1 x_1 + \dots + \theta_D x_D = \theta^\top x$
- But we could take $h(x) = \theta_1 \phi_1(x) + \cdots + \theta_D \phi_D(x) = \theta^\top \phi(x)$
 - Now x can be any dimension (and need not be real-valued), and D is the dimension of the features $\phi(x)$

- Linear models can be very flexible given non-trivial features
- We've been considering $h(x) = \theta_1 x_1 + \dots + \theta_D x_D = \theta^\top x$
- But we could take $h(x) = \theta_1 \phi_1(x) + \cdots + \theta_D \phi_D(x) = \theta^\top \phi(x)$
 - Now x can be any dimension (and need not be real-valued), and D is the dimension of the features $\phi(x)$
 - E.g. for $x \in \mathbb{R}$, $\phi(x) = [1, x, x^2]^\top$

- Linear models can be very flexible given non-trivial features
- We've been considering $h(x) = \theta_1 x_1 + \dots + \theta_D x_D = \theta^\top x$
- But we could take $h(x) = \theta_1 \phi_1(x) + \cdots + \theta_D \phi_D(x) = \theta^\top \phi(x)$
 - Now x can be any dimension (and need not be real-valued), and D is the dimension of the features $\phi(x)$
 - E.g. for $x \in \mathbb{R}$, $\phi(x) = [1, x, x^2]^{\top}$ recall demo

- Linear models can be very flexible given non-trivial features
- We've been considering $h(x) = \theta_1 x_1 + \dots + \theta_D x_D = \theta^\top x$
- But we could take $h(x) = \theta_1 \phi_1(x) + \cdots + \theta_D \phi_D(x) = \theta^\top \phi(x)$
 - Now x can be any dimension (and need not be real-valued), and D is the dimension of the features $\phi(x)$
 - E.g. for $x \in \mathbb{R}$, $\phi(x) = [1, x, x^2]^{\top}$ recall demo
 - More generally, for $x \in \mathbb{R}^{D_x}$, $\phi(x)$ could collect all polynomials of some degree r or smaller

- Linear models can be very flexible given non-trivial features
- We've been considering $h(x) = \theta_1 x_1 + \dots + \theta_D x_D = \theta^\top x$
- But we could take $h(x) = \theta_1 \phi_1(x) + \cdots + \theta_D \phi_D(x) = \theta^\top \phi(x)$
 - Now x can be any dimension (and need not be real-valued), and D is the dimension of the features $\phi(x)$
 - E.g. for $x \in \mathbb{R}$, $\phi(x) = [1, x, x^2]^{\top}$ recall demo
 - More generally, for $x \in \mathbb{R}^{D_x}$, $\phi(x)$ could collect all polynomials of some degree r or smaller
 - Lots of other options: e.g. Fourier basis, logistic basis, wavelets, splines

- Linear models can be very flexible given non-trivial features
- We've been considering $h(x) = \theta_1 x_1 + \dots + \theta_D x_D = \theta^\top x$
- But we could take $h(x) = \theta_1 \phi_1(x) + \cdots + \theta_D \phi_D(x) = \theta^\top \phi(x)$
 - Now x can be any dimension (and need not be real-valued), and D is the dimension of the features $\phi(x)$
 - E.g. for $x \in \mathbb{R}$, $\phi(x) = [1, x, x^2]^{\top}$ recall demo
 - More generally, for $x \in \mathbb{R}^{D_x}$, $\phi(x)$ could collect all polynomials of some degree r or smaller
 - Lots of other options: e.g. Fourier basis, logistic basis, wavelets, splines
 - For any fixed $\phi(x)$, basically all the math we did goes through with $\phi(x)$ in place of x

- Linear models can be very flexible given non-trivial features
- We've been considering $h(x) = \theta_1 x_1 + \dots + \theta_D x_D = \theta^\top x$
- But we could take $h(x) = \theta_1 \phi_1(x) + \cdots + \theta_D \phi_D(x) = \theta^\top \phi(x)$
 - Now x can be any dimension (and need not be real-valued), and D is the dimension of the features $\phi(x)$
 - E.g. for $x \in \mathbb{R}$, $\phi(x) = [1, x, x^2]^{\top}$ recall demo
 - More generally, for $x \in \mathbb{R}^{D_x}$, $\phi(x)$ could collect all polynomials of some degree r or smaller
 - Lots of other options: e.g. Fourier basis, logistic basis, wavelets, splines
 - For any fixed $\phi(x)$, basically all the math we did goes through with $\phi(x)$ in place of x
 - E.g. for NxD matrix Φ with nth row $\phi(x^{(n)})^{\top}$

- Linear models can be very flexible given non-trivial features
- We've been considering $h(x) = \theta_1 x_1 + \cdots + \theta_D x_D = \theta^{\top} x$
- But we could take $h(x) = \theta_1 \phi_1(x) + \cdots + \theta_D \phi_D(x) = \theta^{\top} \phi(x)$
 - Now x can be any dimension (and need not be realvalued), and D is the dimension of the features $\phi(x)$
 - E.g. for $x \in \mathbb{R}$, $\phi(x) = [1, x, x^2]^{\top}$ recall demo
 - More generally, for $x \in \mathbb{R}^{D_x}$, $\phi(x)$ could collect all polynomials of some degree r or smaller
 - Lots of other options: e.g. Fourier basis, logistic basis, wavelets, splines
 - For any fixed $\phi(x)$, basically all the math we did goes through with $\phi(x)$ in place of x
 - E.g. for NxD matrix Φ with nth row $\phi(x^{(n)})^{\top}$ Then OLS: $\hat{\theta} = (\Phi^{\top}\Phi)^{-1}\Phi^{\top}Y$

- Linear models can be very flexible given non-trivial features
- We've been considering $h(x) = \theta_1 x_1 + \dots + \theta_D x_D = \theta^\top x$
- But we could take $h(x) = \theta_1 \phi_1(x) + \cdots + \theta_D \phi_D(x) = \theta^\top \phi(x)$
 - Now x can be any dimension (and need not be real-valued), and D is the dimension of the features $\phi(x)$
 - E.g. for $x \in \mathbb{R}$, $\phi(x) = [1, x, x^2]^{\top}$ recall demo
 - More generally, for $x \in \mathbb{R}^{D_x}$, $\phi(x)$ could collect all polynomials of some degree r or smaller
 - Lots of other options: e.g. Fourier basis, logistic basis, wavelets, splines
 - For any fixed $\phi(x)$, basically all the math we did goes through with $\phi(x)$ in place of x
 - E.g. for NxD matrix Φ with nth row $\phi(x^{(n)})^{\top}$
 - Then OLS: $\hat{\theta} = (\Phi^{\top}\Phi)^{-1}\Phi^{\top}Y$
- Deep neural networks perform nonlinear feature extraction/ learning