

Inductive bias and PCA

COMS 4771 Fall 2019

Overview

- ▶ Inductive biases and regularization
- ▶ Model averaging and Bayesian perspectives
- ▶ Principal component analysis
- ▶ Gradient descent

0 / 39

1 / 39

Inductive bias

- ▶ What if ERM solution is not unique?
- ▶ Infinitely-many solutions to normal equations.
- ▶ Which one should we pick?
 - ▶ Possible answer: Pick shortest solution, i.e., of minimum (squared) Euclidean norm $\|w\|_2^2$.
 - ▶ Smaller norm \Rightarrow slower variations (Cauchy-Schwarz):
$$|w^\top x - w^\top x'| \leq \|w\|_2 \cdot \|x - x'\|_2$$
 - ▶ But data does not give reason to choose shorter w over longer w .
 - ▶ Preference for short w is an example of an [inductive bias](#).
- ▶ All learning algorithms encode some form of inductive bias.

2 / 39

Example of minimum norm inductive bias I

- ▶ Trigonometric feature expansion with particular weighting

$$\varphi(x) = (1, \sin(x), \cos(x), \frac{1}{2} \sin(2x), \frac{1}{2} \cos(2x), \frac{1}{3} \sin(3x), \frac{1}{3} \cos(3x), \dots)$$

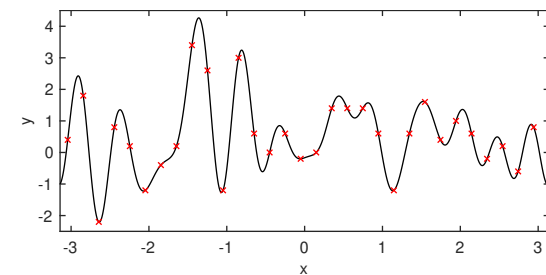


Figure 1: Arbitrary solutions to normal equations can be arbitrarily “wiggly”

3 / 39

Example of minimum norm inductive bias II

- ▶ Trigonometric feature expansion with particular weighting

$$\varphi(x) = (1, \sin(x), \cos(x), \frac{1}{2} \sin(2x), \frac{1}{2} \cos(2x), \frac{1}{3} \sin(3x), \frac{1}{3} \cos(3x), \dots)$$

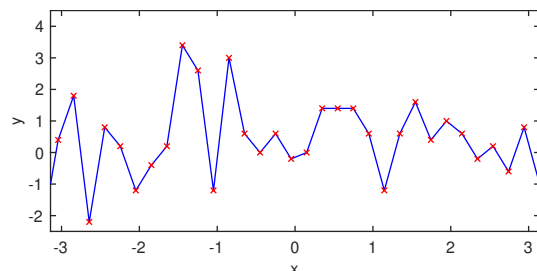


Figure 2: Least norm solution is a very particular interpolation

4 / 39

Representation of minimum norm solution

- ▶ **Claim:** The minimum (Euclidean) norm solution to normal equations lives in span of the x_i 's (i.e., in $\text{range}(\mathbf{A}^\top)$).
 - ▶ I.e., can write

$$\mathbf{w} = \mathbf{A}^\top \boldsymbol{\alpha} = \sum_{i=1}^n \alpha_i \mathbf{x}_i$$

for some $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_n) \in \mathbb{R}^n$.

- ▶ In fact, the solution in $\text{range}(\mathbf{A}^\top)$ is unique!

5 / 39

Regularized ERM

- ▶ Combine two concerns: making both $\widehat{\mathcal{R}}(\mathbf{w})$ and $\|\mathbf{w}\|_2^2$ small
 - ▶ Pick $\lambda \geq 0$, and minimize $\widehat{\mathcal{R}}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2$
- ▶ If $\lambda > 0$, solution is always unique (even if $n < d$).
 - ▶ Called [ridge regression](#).
 - ▶ $\lambda = 0$ is ERM/OLS.
 - ▶ λ controls how much to pay attention to [regularizer](#) $\|\mathbf{w}\|_2^2$ relative to [data fitting term](#) $\widehat{\mathcal{R}}(\mathbf{w})$
 - ▶ λ is hyperparameter to tune (e.g., using cross-validation)

6 / 39

Data augmentation I

- ▶ Let $\tilde{\mathbf{A}} = \begin{bmatrix} \mathbf{A} \\ \sqrt{\lambda} \mathbf{I} \end{bmatrix}$ and $\tilde{\mathbf{b}} = \begin{bmatrix} \mathbf{b} \\ \mathbf{0} \end{bmatrix} \in \mathbb{R}^{n+d}$
- ▶ Then $\|\tilde{\mathbf{A}}\mathbf{w} - \tilde{\mathbf{b}}\|_2^2 = \widehat{\mathcal{R}}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2$ (ridge regression objective)
- ▶ Interpretation:
 - ▶ d "fake" data points, ensures augmented data matrix $\tilde{\mathbf{A}}$ has rank d
 - ▶ All corresponding labels are zero.

- ▶ So ridge regression solution is $\hat{\mathbf{w}} =$

7 / 39

- ▶ Domain-specific data augmentation: e.g., image transformations

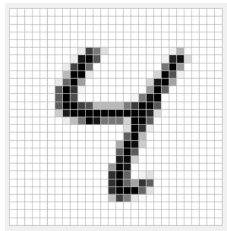
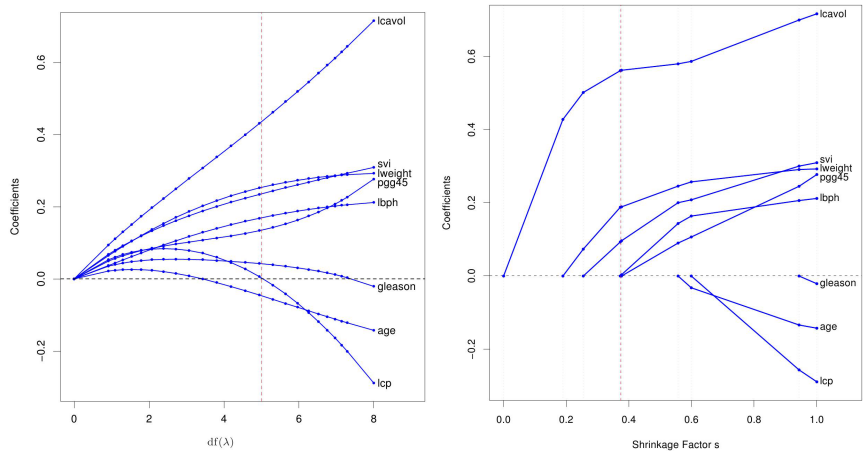


Figure 3: Pixels of OCR image

- ▶ Lasso: minimize $\hat{\mathcal{R}}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1$
 - ▶ Here, $\|\mathbf{v}\|_1 = \sum_{i=1}^n |v_i|$, sum of absolute values of vector entries
 - ▶ Prefers short \mathbf{w} , where length is measured using different norm
 - ▶ Tends to produce \mathbf{w} that are sparse (i.e., have few non-zero entries), or at least are well-approximated by sparse vectors.
 - ▶ A different inductive bias

- ▶ Example: coefficient profile of Lasso vs ridge
- ▶ \mathbf{X} = clinical measurements, Y = level of prostate cancer antigen
- ▶ Horizontal axis: varying λ (large λ to left, small λ to right).
- ▶ Vertical axis: coefficient value in ridge and Lasso solutions, for eight different features



- ▶ **Theorem:** Pick any $\mathbf{w} \in \mathbb{R}^d$ and any $\varepsilon \in (0, 1)$. Form $\tilde{\mathbf{w}} \in \mathbb{R}^d$ by including the $\lceil 1/\varepsilon^2 \rceil$ largest (by magnitude) coefficients of \mathbf{w} , and setting remaining entries to zero. Then

$$\|\tilde{\mathbf{w}} - \mathbf{w}\|_2 \leq \varepsilon \|\mathbf{w}\|_1.$$

- ▶ If $\|\mathbf{w}\|_1$ is small (compared to $\|\mathbf{w}\|_2$), then theorem says \mathbf{w} is well-approximated by sparse vector.

- ▶ Lasso also tries to make coefficients small. What if we only care about sparsity?
- ▶ Subset selection: minimize empirical risk among all k -sparse solutions
- ▶ Greedy algorithms: repeatedly choose new variables to “include” in support of w until k variables are included.
 - ▶ Forward stepwise regression / orthogonal matching pursuit: Each time you “include” a new variable, re-fit all coefficients for included variables.
 - ▶ Often works as well as Lasso
- ▶ Why do we care about sparsity?

- ▶ Suppose we have M real-valued predictors, $\hat{f}_1, \dots, \hat{f}_M$
 - ▶ E.g., nearest neighbor regression, regression trees, linear models with different feature expansions, ...
- ▶ How to take advantage of all of them?
- ▶ Model selection: pick the best one, e.g., using hold-out method or K -fold cross-validation
- ▶ Model averaging: form “ensemble” predictor \hat{f}_{avg} , where for any x ,

$$\hat{f}_{\text{avg}}(x) := \frac{1}{M} \sum_{i=1}^M \hat{f}_i(x).$$

- ▶ **Theorem**: Risk of \hat{f}_{avg} :

$$\mathcal{R}(\hat{f}_{\text{avg}}) = \frac{1}{M} \sum_{i=1}^M \mathcal{R}(\hat{f}_i) - \frac{1}{M} \sum_{i=1}^M \mathbb{E} \left[(\hat{f}_{\text{avg}}(X) - \hat{f}_i(X))^2 \right].$$

- ▶ Better than model selection when:
 - ▶ all \hat{f}_i have similar risks, and
 - ▶ all \hat{f}_i predict very differently from each other

- ▶ In model averaging, “weights” of $1/M$ for all \hat{f}_i seems arbitrary
- ▶ Can “learn” weights using linear regression!
 - ▶ Use feature expansion $\varphi(x) = (\hat{f}_1(x), \dots, \hat{f}_M(x))$
 - ▶ Called stacking
 - ▶ Use additional data (independent of $\hat{f}_1, \dots, \hat{f}_M$)
- ▶ Upshot: Any function (even learned functions) can be a feature
- ▶ Conversely: Behind every feature is a deliberate modeling choice

- Bayesian inference: probabilistic approach to updating beliefs
 - Posit a (parametric) statistical model for data (likelihood)
 - Start with some beliefs about the parameters of model (prior)
 - Update beliefs after seeing data (posterior)

$$\underbrace{\Pr(\mathbf{w} \mid \text{data})}_{\text{posterior}(\mathbf{w})} \propto \underbrace{\Pr(\mathbf{w})}_{\text{prior}(\mathbf{w})} \cdot \underbrace{\Pr(\text{data} \mid \mathbf{w})}_{\text{likelihood}(\mathbf{w})}$$

- (Finding proportionality constant is often the computationally challenging part of belief updating.)
- Basis for reasoning in humans (maybe?), robots, etc.

16 / 39

- Can use Bayesian inference framework for designing estimation/learning algorithms (even if you aren't a Bayesian!)
 - E.g., instead of computing entire posterior distribution, find the \mathbf{w} with highest posterior probability
 - Called maximum a posteriori (MAP) estimator
 - Just find \mathbf{w} to maximize

$$\text{prior}(\mathbf{w}) \times \text{likelihood}(\mathbf{w}).$$

- (Avoids issue with finding proportionality constant.)

17 / 39

- In linear regression model, express prior belief about $\mathbf{w} = (w_1, \dots, w_d)$ using a probability distribution with density function π
 - Simple choice: $\text{prior}(w_1, \dots, w_d) = \prod_{j=1}^d \sqrt{\frac{\tau}{2\pi}} \exp(-\tau w_j^2/2)$
 - I.e., treat w_1, \dots, w_d as independent $N(0, 1/\tau)$ random variables
- Likelihood model: $(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)$ are conditionally independent given \mathbf{w} , and $Y_i \mid (\mathbf{X}_i, \mathbf{w}) \sim N(\mathbf{X}_i^\top \mathbf{w}, 1)$.
- What is the MAP?

18 / 39

- Find \mathbf{w} to maximize

$$\underbrace{\prod_{j=1}^d \sqrt{\frac{\tau}{2\pi}} \exp(-\tau w_j^2/2)}_{\text{prior}(\mathbf{w})} \cdot \underbrace{\prod_{i=1}^n p(x_i) \cdot \frac{1}{\sqrt{2\pi}} \exp(-(y_i - \mathbf{x}_i^\top \mathbf{w})^2/2)}_{\text{likelihood}(\mathbf{w})}.$$

(Here, p is marginal density of \mathbf{X} ; unimportant.)

- Take logarithm and omit terms not involving \mathbf{w} :

$$-\frac{\tau}{2} \sum_{j=1}^d w_j^2 - \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \mathbf{w})^2.$$

- For $\tau = n\lambda$, same as minimizing

$$\frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^\top \mathbf{w} - y_i)^2 + \lambda \|\mathbf{w}\|_2^2,$$

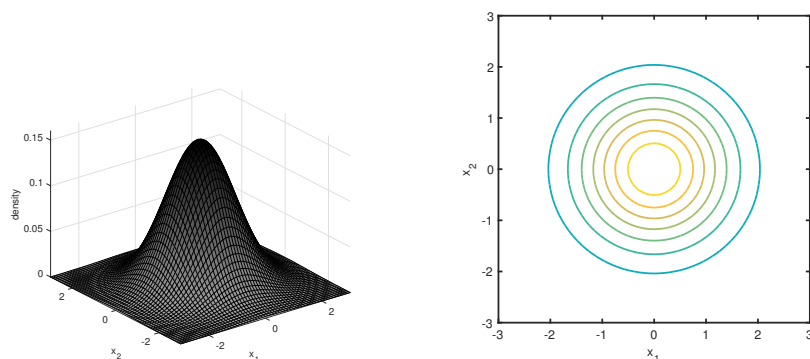
which is the ridge regression objective!

- What about different Gaussian prior?

19 / 39

Multivariate Gaussians I: Isotropic Gaussians

- Start with $\mathbf{X} = (X_1, \dots, X_d) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, i.e., X_1, \dots, X_d are iid $\mathcal{N}(0, 1)$ random variables.
 - Probability density function is product of (univariate) Gaussian densities
 - $\mathbb{E}(X_i) = 0$
 - $\text{var}(X_i) = \text{cov}(X_i, X_i) = 1$, $\text{cov}(X_i, X_j) = 0$ for $i \neq j$
 - Arrange in mean vector $\mathbb{E}(\mathbf{X}) = \mathbf{0}$, covariance matrix $\text{cov}(\mathbf{X}) = \mathbf{I}$



20 / 39

Affine transformations of random vectors

- Start with any random vector \mathbf{X} , then apply linear transformation, followed by translation
- $\mathbf{Y} := \mathbf{M}\mathbf{X} + \boldsymbol{\mu}$, for $\mathbf{M} \in \mathbb{R}^{k \times d}$ and $\boldsymbol{\mu} \in \mathbb{R}^k$
- $\mathbb{E}(\mathbf{Y}) =$
- $\text{cov}(\mathbf{Y}) =$
- Let $\mathbf{u} \in \mathbb{R}^d$ be a unit vector ($\|\mathbf{u}\|_2 = 1$), and $\mathbf{Y} := \mathbf{u}^\top \mathbf{X}$ (projection of \mathbf{X} along direction \mathbf{u}).
- $\mathbb{E}(\mathbf{Y}) =$
- $\text{var}(\mathbf{Y}) =$

21 / 39

Multivariate Gaussians II: General Gaussians

- If $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $\mathbf{Y} = \mathbf{M}\mathbf{X} + \boldsymbol{\mu}$, we have $\mathbb{E}(\mathbf{Y}) = \boldsymbol{\mu}$ and $\text{cov}(\mathbf{Y}) = \mathbf{M}\mathbf{M}^\top$
 - Assume $\mathbf{M} \in \mathbb{R}^{d \times d}$ is invertible (else we get a degenerate Gaussian distribution).
 - We say $\mathbf{Y} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{M}\mathbf{M}^\top)$
 - Density function given by

$$\frac{1}{(2\pi)^{d/2} |\mathbf{M}\mathbf{M}^\top|^{1/2}} \exp\left(-\frac{1}{2} \|\mathbf{M}^{-1}(\mathbf{y} - \boldsymbol{\mu})\|_2^2\right).$$

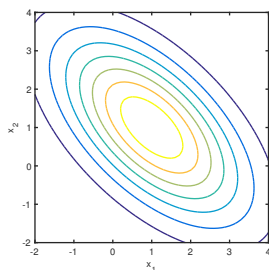


Figure 4: Contour lines of a bivariate Gaussian density

22 / 39

MAP with general Gaussian priors

- Prior: multivariate Gaussian, written $\mathbf{w} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
 - Probability density is $\text{prior}(\mathbf{w}) \propto \exp\left(-\frac{1}{2}(\mathbf{w} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{w} - \boldsymbol{\mu})\right)$
- Find \mathbf{w} to maximize

$$\underbrace{\exp\left(-\frac{1}{2}(\mathbf{w} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{w} - \boldsymbol{\mu})\right)}_{\text{prior}(\mathbf{w})} \cdot \underbrace{\prod_{i=1}^n p(\mathbf{x}_i) \cdot \frac{1}{\sqrt{2\pi}} \exp(-(y_i - \mathbf{x}_i^\top \mathbf{w})^2/2)}_{\text{likelihood}(\mathbf{w})}.$$
- Take logarithm and omit terms not involving \mathbf{w} :

$$-\frac{1}{2}(\mathbf{w} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{w} - \boldsymbol{\mu}) - \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \mathbf{w})^2.$$
- For $\mathbf{C} := \boldsymbol{\Sigma}^{-1}/n$, same as minimizing

$$\frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^\top \mathbf{w} - y_i)^2 + (\mathbf{w} - \boldsymbol{\mu})^\top \mathbf{C}(\mathbf{w} - \boldsymbol{\mu}),$$
 a different regularizer!

23 / 39

Eigendecomposition

- ▶ Every symmetric matrix $M \in \mathbb{R}^{d \times d}$ has d real eigenvalues, which we arrange as $\lambda_1 \geq \dots \geq \lambda_d$
 - ▶ Can choose corresponding eigenvectors $v_1, \dots, v_d \in \mathbb{R}^d$ to be orthonormal
 - ▶ This means $Mv_i = \lambda_i v_i$ for each $i = 1, \dots, d$, and $v_i^T v_j = \mathbb{1}_{\{i=j\}}$
- ▶ Often arrange v_1, \dots, v_d in an orthogonal matrix $V := [v_1 | \dots | v_d]$
 - ▶ $V^T V = I$ and $V V^T = \sum_{i=1}^d v_i v_i^T = I$
- ▶ Eigendecomposition (spectral decomposition):

▶ Diagonalization:

24 / 39

Covariance matrix

- ▶ $A \in \mathbb{R}^{n \times d}$ is data matrix
- ▶ $\Sigma := A^T A = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$ is (empirical) second-moment matrix
 - ▶ Assuming $\frac{1}{n} \sum_{i=1}^n x_i = 0$, this is the (empirical) covariance matrix
- ▶ For any unit vector $u \in \mathbb{R}^d$,

$$u^T \Sigma u = \frac{1}{n} \sum_{i=1}^n (u^T x_i)^2$$

is variance of data along direction u

- ▶ Note: some pixels in OCR data have very little (or zero!) variation

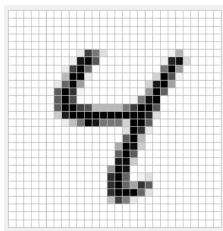


Figure 5: Pixels of OCR image

25 / 39

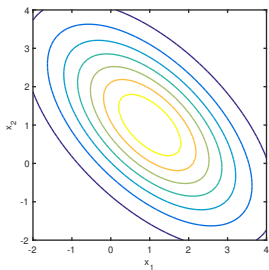
Top eigenvector

- ▶ Σ is symmetric, so can write eigendecomposition

$$\Sigma = \sum_{i=1}^n \lambda_i v_i v_i^T$$

- ▶ In which direction is variance maximized?
- ▶ Answer: v_1 , corresponding to largest eigenvalue λ_1
 - ▶ Called the top eigenvector
 - ▶ This follows from the following characterization of v_1 :

$$v_1^T \Sigma v_1 = \max_{u \in \mathbb{R}^d: \|u\|_2=1} u^T \Sigma u = \lambda_1.$$



26 / 39

Top k eigenvectors

- ▶ What about among directions orthogonal to v_1 ?
 - ▶ Answer: v_2 , corresponding to second largest eigenvalue λ_2
- ▶ Etc.
- ▶ For any k , $V_k := [v_1 | \dots | v_k]$ satisfies

$$\sum_{i=1}^k v_i^T \Sigma v_i = \text{tr}(V_k^T \Sigma V_k) = \max_{U \in \mathbb{R}^{d \times k}: U^T U = I} \text{tr}(U^T \Sigma U) = \sum_{i=1}^k \lambda_i$$

(the top k eigenvectors)

27 / 39

Principal component analysis

- k -dimensional principal components analysis (PCA) mapping:

$$\varphi(\mathbf{x}) = (\mathbf{x}^\top \mathbf{v}_1, \dots, \mathbf{x}^\top \mathbf{v}_k) = \mathbf{V}_k^\top \mathbf{x} \in \mathbb{R}^k$$

where $\mathbf{V}_k = [\mathbf{v}_1 | \dots | \mathbf{v}_k] \in \mathbb{R}^{d \times k}$

- (Only really makes sense when $\lambda_k > 0$.)
- This is a form of dimensionality reduction when $k < d$.

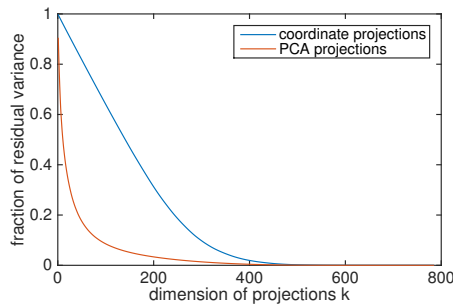


Figure 7: Fraction of residual variance with PCA and coordinate projections

28 / 39

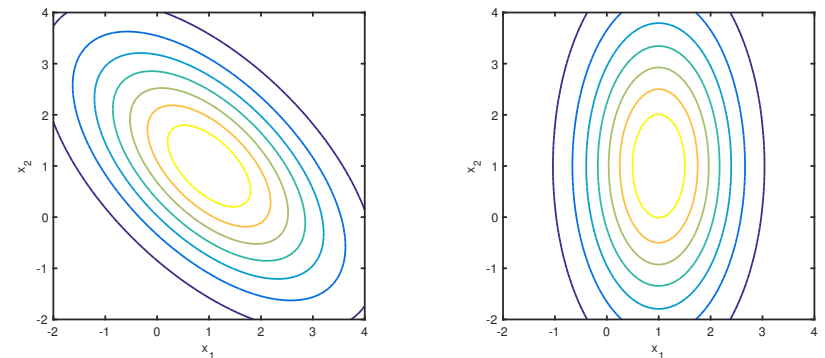
Covariance of data upon PCA mapping

- Covariance of data upon PCA mapping:

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \varphi(\mathbf{x}_i) \varphi(\mathbf{x}_i)^\top &= \frac{1}{n} \sum_{i=1}^n \mathbf{V}_k^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{V}_k \\ &= \mathbf{V}_k^\top \Sigma \mathbf{V}_k \\ &= \mathbf{\Lambda}_k \end{aligned}$$

where $\mathbf{\Lambda}_k$ is diagonal matrix with $\lambda_1, \dots, \lambda_k$ along diagonal.

- In particular, coordinates in $\varphi(\mathbf{x})$ -representation are uncorrelated.



29 / 39

PCA and linear regression

- Use k -dimensional PCA mapping $\varphi(\mathbf{x}) = \mathbf{V}_k^\top \mathbf{x}$ with OLS
- (Assume rank of \mathbf{A} is at least k , so $\mathbf{A}^\top \mathbf{A}$ has $\lambda_k > 0$)
- Data matrix is

$$\frac{1}{\sqrt{n}} \begin{bmatrix} \leftarrow & \varphi(\mathbf{x}_1)^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \varphi(\mathbf{x}_n)^\top & \rightarrow \end{bmatrix} = \frac{1}{\sqrt{n}} \begin{bmatrix} \leftarrow & \mathbf{x}_1^\top \mathbf{V}_k & \rightarrow \\ & \vdots & \\ \leftarrow & \mathbf{x}_n^\top \mathbf{V}_k & \rightarrow \end{bmatrix} = \mathbf{A} \mathbf{V}_k \in \mathbb{R}^{n \times k}$$

- Therefore, OLS solution is

$$\begin{aligned} \hat{\beta} &= (\mathbf{V}_k^\top \mathbf{A}^\top \mathbf{A} \mathbf{V}_k)^{-1} (\mathbf{A} \mathbf{V}_k)^\top \mathbf{b} \\ &= \mathbf{\Lambda}_k^{-1} \mathbf{V}_k^\top \mathbf{A}^\top \mathbf{b} \end{aligned}$$

(Note: here $\hat{\beta} \in \mathbb{R}^k$.)

30 / 39

Principal component regression

- Use $\hat{\beta} = \mathbf{\Lambda}_k^{-1} \mathbf{V}_k^\top \mathbf{A}^\top \mathbf{b}$ to predict on new $\mathbf{x} \in \mathbb{R}^d$:

$$\begin{aligned} \varphi(\mathbf{x})^\top \hat{\beta} &= (\mathbf{V}_k^\top \mathbf{x})^\top \mathbf{\Lambda}_k^{-1} \mathbf{V}_k^\top \mathbf{A}^\top \mathbf{b} \\ &= \mathbf{x}^\top (\mathbf{V}_k \mathbf{\Lambda}_k^{-1} \mathbf{V}_k^\top) (\mathbf{A}^\top \mathbf{b}) \end{aligned}$$

- So “effective” weight vector (that acts directly on \mathbf{x} rather than $\varphi(\mathbf{x})$) is given by

$$\hat{\mathbf{w}} := (\mathbf{V}_k \mathbf{\Lambda}_k^{-1} \mathbf{V}_k^\top) (\mathbf{A}^\top \mathbf{b}) = \left(\sum_{i=1}^k \frac{1}{\lambda_i} \mathbf{v}_i \mathbf{v}_i^\top \right) (\mathbf{A}^\top \mathbf{b}).$$

- This is called principal component regression (PCR) (here, k is hyperparameter)
- Alternative hyper-parameterization: $\lambda > 0$; same as before but using the largest k such that $\lambda_k \geq \lambda$.

31 / 39

Spectral regularization

- PCR and ridge regression are examples of [spectral regularization](#).
- For a function $g: \mathbb{R} \rightarrow \mathbb{R}$, write $g(\mathbf{M})$ to mean

$$g(\mathbf{M}) = \sum_{i=1}^d g(\lambda_i) \mathbf{v}_i \mathbf{v}_i^\top$$

where \mathbf{M} has eigendecomposition $\mathbf{M} = \sum_{i=1}^d \lambda_i \mathbf{v}_i \mathbf{v}_i^\top$.

- **Claim:** Can write each of PCR and ridge regression as

$$\hat{\mathbf{w}} = g(\mathbf{A}^\top \mathbf{A}) \mathbf{A}^\top \mathbf{b}$$

for appropriate function g (depending on λ).

32 / 39

Comparing ridge regression and PCR

Ridge: $g(z) = \frac{1}{z + \lambda}$; PCR: $g(z) = \mathbb{1}_{\{z \geq \lambda\}} \cdot \frac{1}{z}$

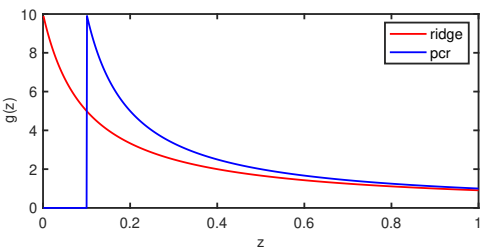


Figure 8: Spectral regularization function g for ridge and PCR ($\lambda = 0.1$)

- Interpretation:
 - PCR only uses directions with sufficient variability; ignores the rest
 - Ridge artificially inflates the variance in all directions

33 / 39

Optimization for linear regression

- Back to considering ordinary least squares.
- Gaussian elimination to solve normal equations can be slow when d is large (time is $O(nd^2)$).
- Alternative: find approximate solution using [gradient descent](#)

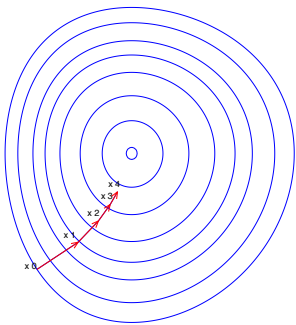


Figure 9: Gradient descent

34 / 39

Gradient descent for linear regression

- Algorithm: start with some $\mathbf{w}^{(0)} \in \mathbb{R}^d$ and $\eta > 0$.
 - For $t = 1, 2, \dots$:

$$\begin{aligned} \mathbf{w}^{(t)} &:= \mathbf{w}^{(t-1)} - 2\eta \mathbf{A}^\top (\mathbf{A} \mathbf{w}^{(t-1)} - \mathbf{b}) \\ &= \mathbf{w}^{(t-1)} - 2\eta \cdot \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^\top \mathbf{w}^{(t-1)} - y_i) \mathbf{x}_i \end{aligned}$$

- Time to multiply matrix by vector is linear in matrix size.
- So each iteration takes time $O(nd)$.
- η is called [step size](#) (somewhat of a misnomer)

35 / 39

Motivation for gradient descent

- Why move in direction of (negative) gradient?
- Affine approximation of $\widehat{\mathcal{R}}(\mathbf{w} + \delta)$ around \mathbf{w} :

- Use $\delta := -\eta \nabla \widehat{\mathcal{R}}(\mathbf{w})$ for some $\eta > 0$:

36 / 39

Interpretation of gradient descent for linear regression

- Interpretation (specific to least squares objective):

$$\nabla (\mathbf{x}_i^\top \mathbf{w} - y_i)^2 = 2(\mathbf{x}_i^\top \mathbf{w} - y_i) \mathbf{x}_i.$$

- If $\mathbf{x}_i^\top \mathbf{w} > y_i$, subtract a little bit of \mathbf{x}_i from \mathbf{w}
- If $\mathbf{x}_i^\top \mathbf{w} < y_i$, add a little bit of \mathbf{x}_i from \mathbf{w}
- If $\mathbf{x}_i^\top \mathbf{w} = y_i$, i -th term has no contribution

37 / 39

Behavior of gradient descent for linear regression

- **Theorem:** Let $\hat{\mathbf{w}}$ be the minimum Euclidean norm solution to normal equations. Assume $\mathbf{w}^{(0)} = \mathbf{0}$. Write eigendecomposition $\mathbf{A}^\top \mathbf{A} = \sum_{i=1}^r \lambda_i \mathbf{v}_i \mathbf{v}_i^\top$ with $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0$. Then $\mathbf{w}^{(t)} \in \text{range}(\mathbf{A}^\top)$ and

$$\mathbf{v}_i^\top \mathbf{w}^{(t)} = \left(2\eta \lambda_i \sum_{k=0}^{t-1} (1 - 2\eta \lambda_i)^k \right) \mathbf{v}_i^\top \hat{\mathbf{w}}, \quad i = 1, \dots, r.$$

- Implications:

- If we choose η such that $2\eta \lambda_i < 1$, then

$$2\eta \lambda_i \sum_{k=0}^{t-1} (1 - 2\eta \lambda_i)^k = 1 - (1 - 2\eta \lambda_i)^t,$$

which converges to 1 as $t \rightarrow \infty$.

- So, when $2\eta \lambda_1 < 1$, we have $\mathbf{w}^{(t)} \rightarrow \hat{\mathbf{w}}$ as $t \rightarrow \infty$.
- Rate of convergence is geometric, i.e., “exponentially fast convergence”.

38 / 39

Inductive bias of gradient descent

- Gradient descent for linear regression has an inductive bias—converges to the minimum norm solution.
- Also a form of spectral regularization, with function

$$g(z) = \mathbb{1}_{\{z > 0\}} \cdot \frac{1 - (1 - 2\eta z)^t}{z}.$$

- Minimum norm solution uses

$$g(z) = \mathbb{1}_{\{z > 0\}} \cdot \frac{1}{z}.$$

39 / 39