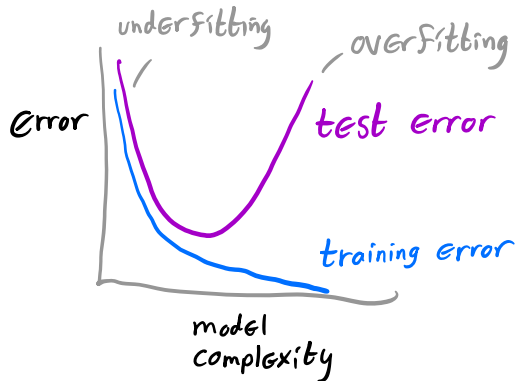**Day 11 - Ridge Regression**

Agenda:

- Review - Bias Variance Tradeoff
- Ridge Regression
- Analytical Formula for Solution to Ridge Regression
- Background - Singular Value Decompositions
- Ridge Regression and Bias Variance Tradeoff

# Bias - Variance Tradeoff

Standard Statistical ML Story:

underfitting     overfitting

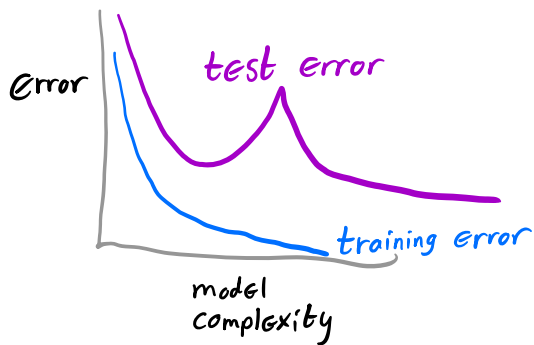Error     test error

training error

model
complexity

higher complexity models
have lower bias but
higher variance

If complexity is too high,
it overfits data, variance term
dominates test error

after a certain threshold,
"larger models are worse"

Modern Story based on Neural Nets:

Error     test error

training error

model
complexity

Underparameterized     Overparameterized
regime                  regime

Test error can decrease as
model complexity continues increasing.

And it can be lower than in
underparameterized regime

Phenomenon: double descent

"larger models are better"

## Ridge Regression

So far, we have used MLE to estimate model
parameters from data

Concern: Overfitting

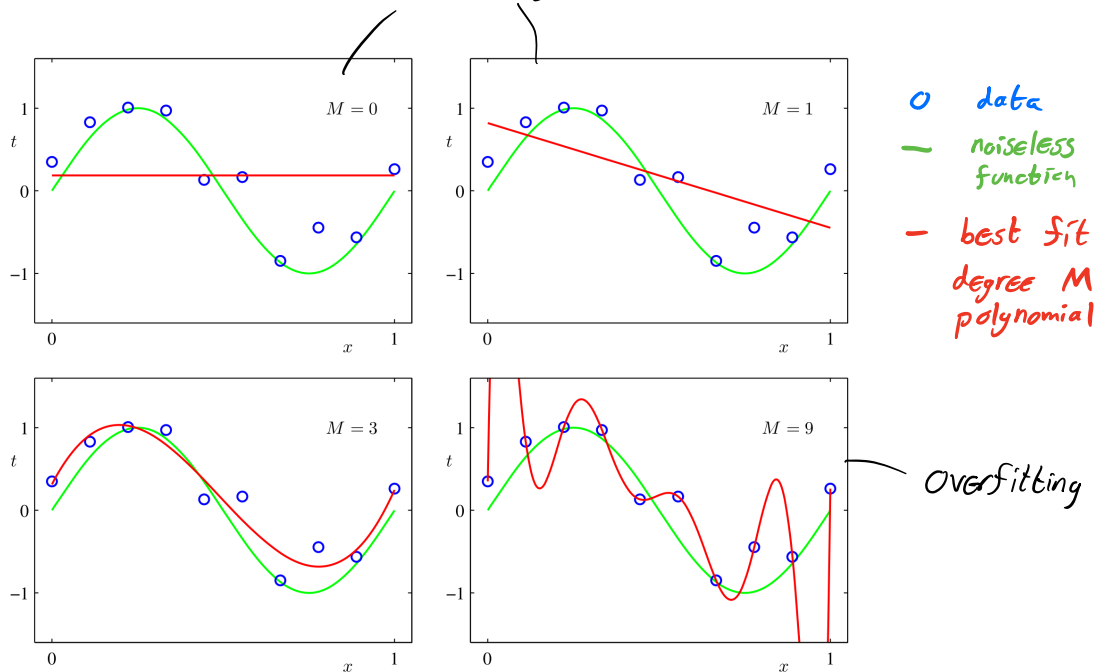Example: Fitting data w/ a degree M polynomial

underfitting



O  data

—  noiseless
   function

—  best fit
   degree M
   polynomial

Overfitting

**Figure 1.4** Plots of polynomials having various orders $M$, shown as red curves, fitted to the data set shown in Figure 1.2.

One way to reduce overfitting,
    use a hypothesis class with lower complexity
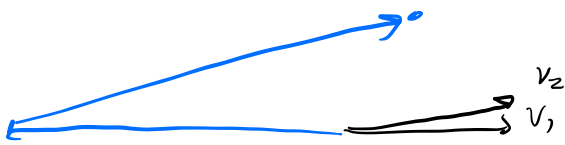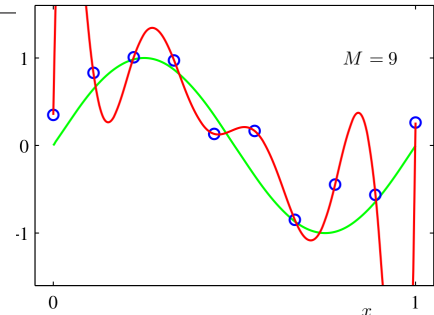    (fewer unknown parameters)

Another way,
    add regularization

A possible indication of overfitting
is having very large learned parameters.
   This often happens when features are highly correlated

**Table 1.2** Table of the coefficients $\mathbf{w}^\star$ for $M = 9$ polynomials with various values for the regularization parameter $\lambda$. Note that $\ln \lambda = -\infty$ corresponds to a model with no regularization, i.e., to the graph at the bottom right in Figure 1.4. We see that, as the value of $\lambda$ increases, the typical magnitude of the coefficients gets smaller.

| | $\ln \lambda = -\infty$ |
|---|---|
| $w_0^\star$ | 0.35 |
| $w_1^\star$ | 232.37 |
| $w_2^\star$ | -5321.83 |
| $w_3^\star$ | 48568.31 |
| $w_4^\star$ | -231639.30 |
| $w_5^\star$ | 640042.26 |
| $w_6^\star$ | -1061800.52 |
| $w_7^\star$ | 1042400.18 |
| $w_8^\star$ | -557682.99 |
| $w_9^\star$ | 125201.43 |

$v_2$
$v_1$

Idea: penalize predictors that have large values of unknown parameters

New formulation for least squares:
   Given data $\{(x_i, y_i)\}_{i=1\cdots n}$  w/  $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$
   where  $y = X\theta + \varepsilon$    w/  $\varepsilon \in \mathbb{R}^n$ has $\mathcal{N}(0, \sigma^2)$ entries

   Estimate $\theta$ by solving          ↗ ridge regression problem

$$\min_\theta \quad \| y - X\theta \|^2 + \lambda \|\theta\|^2$$

$\ell_2$ penalization / $\ell_2$ regularization / weight decay
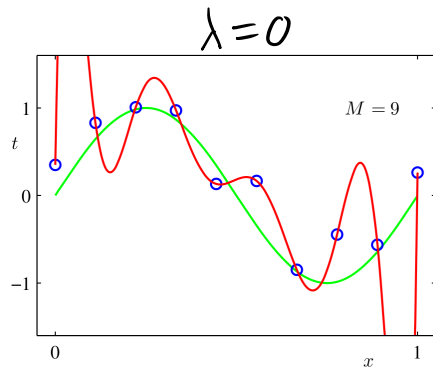
   Solution is given by

$$\hat{\theta}_{ridge} = \left( X^t X + \lambda I_{d \times d} \right)^{-1} X^t y$$

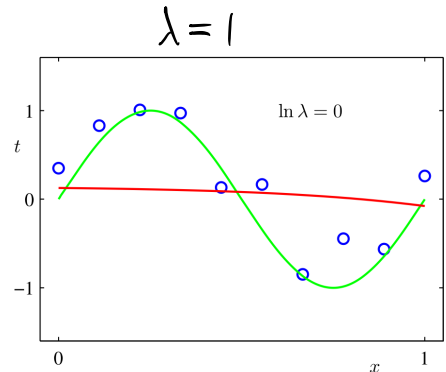w/  $I_{d \times d} = d \times d$ Identity matrix $= \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$
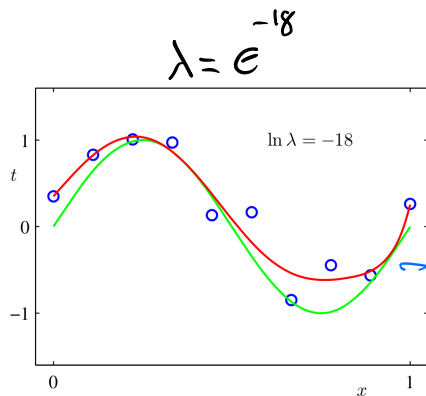
What do solutions look like?

$\lambda = 0$



$M = 9$

$\lambda$ is too low

$\lambda = 1$



$\ln \lambda = 0$

$\lambda$ is too high

$\lambda = e^{-18}$



$\ln \lambda = -18$

$\lambda$ is about right

$$\min_{\theta} \| y - X\theta \|^2 + \lambda \|\theta\|^2$$

degree 9 polynomial

**Solution to ridge regression problem**

Let $y \in \mathbb{R}^n$, $X \in \mathbb{R}^{n \times d}$

The unique solution to

$$\min_{\theta \in \mathbb{R}^d} \| y - X\theta \|^2 + \lambda \|\theta\|^2$$

is

$$\hat{\theta}_{ridge} = \left( X^t X + \lambda I_{d \times d} \right)^{-1} X^t y$$

Proof:    Let $f(\theta) = \| X\theta - y \|^2 + \lambda \| \theta \|^2$

$\nabla f(\theta) = 2 X^t ( X\theta - y) + 2\lambda \theta$

set $\nabla f(\theta) = 0$

$\Rightarrow 2 X^t (X\theta - y) + 2\lambda\theta = 0$

$\Rightarrow X^t X \theta - X^t y + \lambda\theta = 0$

$\Rightarrow (X^t X + \lambda I_{d\times d}) \theta = X^t y$

$\Rightarrow \theta = \underbrace{( X^t X + \lambda I_{d\times d})}^{-1} X^t y.$

Note: this matrix is always invertible if $\lambda > 0$
why?

SVD of a square matrix:

Suppose $A \in \mathbb{R}^{d \times d}$. An SVD of $A$ is given by

$$A = U \Sigma V^t$$

where
- $U$ is $d \times d$ matrix w/ orthonormal columns
- $V$ is $d \times d$ matrix w/ orthonormal columns
- $\Sigma$ is diagonal w/ nonnegative entries $\sigma_1, \sigma_2, \cdots \sigma_d$
  where $\sigma_i \geqslant \sigma_{i+1} \geqslant 0$

The columns of $U$ are the left singular vectors of $A$
— — — $V$ — — right singular vectors — —
The diagonal entries of $\Sigma$ are the singular values of $A$

$$A = \begin{pmatrix} | & | & & | \\ U_1 & U_2 & \cdots & U_d \\ | & | & & | \end{pmatrix} \begin{pmatrix} \sigma_1 & & & O \\ & \sigma_2 & & \\ & & \ddots & \\ O & & & \sigma_d \end{pmatrix} \begin{pmatrix} - & V_1^t & - \\ - & V_2^t & - \\ & \vdots & \\ - & V_d^t & - \end{pmatrix}$$
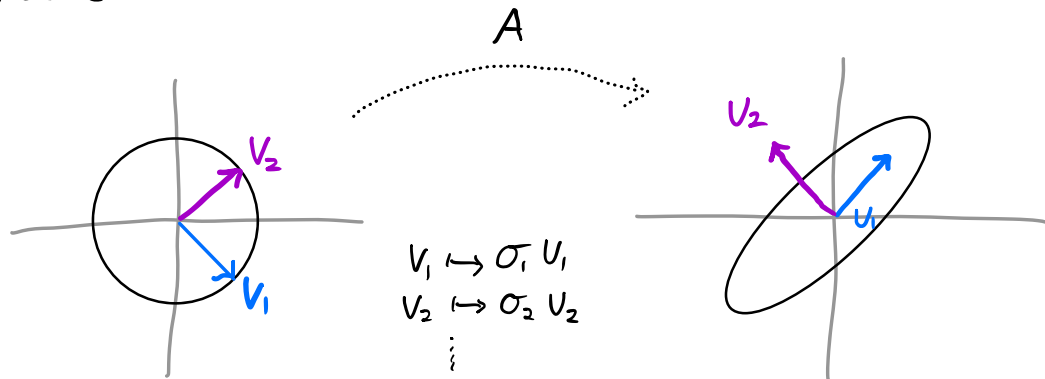
Note: A set $\{ U_1 \cdots U_n \}$ is orthonormal if
- $\| V_i \|^2 = 1$ for all $i$
- $U_i \cdot U_j = 0$ if $i \neq j$

The $ij$ entry of $U^t U = U_i^t U_j = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}$

So $U$ has orthonormal columns if $\underbrace{U^t U}_{d \times d} = I_{d \times d}$

Geometric picture of SVD:

$$A$$



$$V_1 \mapsto \sigma_1 U_1$$
$$V_2 \mapsto \sigma_2 U_2$$
$$\vdots$$

Linear operators map the unit circle to an ellipsoid
The left singular vectors provide the principal axes of
the ellipsoid.

Alternatively, any $A$ is a diagonal matrix if the
domain & range spaces use the right bases.

Given an orthonormal basis $\{V_1 \cdots V_d\}$ of $\mathbb{R}^d$,

if $V = \begin{pmatrix} V_1 & V_2 & \cdots & V_d \end{pmatrix}$ then the coefficients

of $x$ in the basis $\{V_1 \cdots V_d\}$ is given by

$$V^t x. \qquad\qquad V(V^t x) = x$$

So, SVD can be interpreted as

$$A = U \Sigma V^t$$

convert
from basis
given by $U$

diagonal
operator

put input vector
in basis given by $V$

Example

You can use SVD to manipulate matrices easily

Show that if $A \in \mathbb{R}^{n \times n}$ is invertible,

and $A = U \Sigma V^t$ is SVD of $A$, then

$$A^{-1} = V \Sigma^{-1} U^t$$

Proof: If $\Sigma$ is invertible, $\sigma_d > 0$.
Otherwise $V_d$ would be in null space of
$A$, and hence $A$ isn't invertible.

We will show $A(V \Sigma^{-1} U^t) = I_n$.

$$A V \Sigma^{-1} U^t = U \Sigma \underbrace{V^t V}_{I_n} \Sigma^{-1} U^t$$

$$\Sigma = \begin{pmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_d \end{pmatrix}$$

$$\Sigma^{-1} = \begin{pmatrix} \frac{1}{\sigma_1} & & & \\ & \frac{1}{\sigma_2} & & \\ & & \ddots & \\ & & & \frac{1}{\sigma_d} \end{pmatrix}$$

$$= U \underbrace{\Sigma \Sigma^{-1}}_{I_n} V^t$$

$$= U U^t$$

$$= I_n$$

# SVD of a tall rectangular matrix

Let $A \in \mathbb{R}^{n \times d}$ w/ $n \geq d$.

An SVD of $A$ is given by

$$A = U \Sigma V^t$$

w/ $U$ — $n \times d$ matrix w/ orthonormal columns

$\Sigma$ — $d \times d$ diagonal nonnegative matrix
w/ decreasing values along diagonal

$V$ — $d \times d$ matrix w/ orthonormal columns

$$A = \begin{pmatrix} | & & | \\ a_1 & \cdots & a_d \\ | & & | \end{pmatrix} = \begin{pmatrix} | & & | \\ u_1 & \cdots & u_2 \\ | & & | \end{pmatrix} \begin{pmatrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_d \end{pmatrix} \begin{pmatrix} - & v_1^t & - \\ - & v_2^t & - \\ & \vdots & \\ - & v_d^t & - \end{pmatrix}$$

Note $U^t U = I_d$ but $U U^t \neq I_n$ (if $d > n$)

$V^t V = I_d$ & $V V^t = I_d$

**Questions about SVD:**

(a) From an SVD, how can you find the range of a matrix?

(b) From an SVD, how can you find the null space of a matrix?

(c) From an SVD, how can you find the rank of a matrix?

(d) What happens to an SVD if you negate a matrix?

(e) Is the SVD of a matrix unique?

(f) ~~From an SVD of the matrix A, what is an SVD of the matrix A + lambda I?~~

$$\cancel{A = U \Sigma V^t}$$

$$\cancel{A + \lambda I = U \Sigma V^t}$$

(g) What is the relationship of the SVD of a (nonsquare) matrix A with the eigenvector decomposition of $A A^t$ and $A^t A$

$$A = U \Sigma V^t$$

$$A^t A = V \Sigma U^t U \Sigma V^t$$
$$= V \Sigma^2 V^t$$

$$A A^t = U \Sigma V^t V \Sigma^t U^t$$
$$\underbrace{}_{I}$$
$$= U \Sigma \Sigma^t U^t \quad \text{eigenvalue}$$
$$= U \Sigma^2 U^t \quad \text{decomp}$$
$$\text{of } A A^t \quad - \text{eigenvalues}$$
$$\text{are squares of sing values}$$

# Ridge Regression and the Bias Variance Tradeoff

Suppose data $\{(x_i, y_i)\}_{i=1\cdots n}$ follows the distribution

$$y_i = x_i^t \theta^* + \varepsilon_i \quad w/ \quad \varepsilon_i \sim N(0, \sigma^2)$$

That is,

$$y = X\theta^* + \varepsilon$$

Let $X = U\Sigma V^t$ be the SVD of $X$, where $\Sigma = diag(\sigma_1 \cdots \sigma_d)$

The ridge regression estimate of $\theta^*$ is

$$\hat{\theta}_{ridge} = \underbrace{V \, diag\left(\frac{\sigma_1^2}{\sigma_1^2 + \lambda}, \cdots, \frac{\sigma_d^2}{\sigma_d^2 + \lambda}\right) V^t \theta^*}_{Signal \quad \hat{\theta}_{ridge}^{signal}} + \underbrace{V \, diag\left(\frac{\sigma_1}{\sigma_1^2 + \lambda}, \cdots, \frac{\sigma_d}{\sigma_d^2 + \lambda}\right) U^t \varepsilon}_{noise \quad \hat{\theta}_{ridge}^{noise}}.$$

Let's analyze bias and variance of $\hat{\theta}_{ridge}$.

Note: — $\mathbb{E}\,\hat{\theta}_{ridge}^{noise} = 0$. So first term controls bias

— first term doesn't depend on $\varepsilon$. So second term controls Variance

Analyze $\hat{\theta}_{ridge}^{signal}$ — if $\lambda = 0$ $\quad \hat{\theta}_{ridge}^{signal} = VV^t \theta^* = \theta^*$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ Unbiased

$\qquad\qquad\qquad\qquad$ if $\lambda = \infty$ $\quad \hat{\theta}_{ridge}^{signal} = 0 \qquad$ biased

$$\text{Bias increases with } \lambda.$$

Analyze $\hat{\theta}_{ridge}^{noise}$ — if $\lambda = \infty$ $\quad \hat{\theta}_{ridge}^{noise} = 0$ $\quad$ low variance

$$\text{if } \lambda = 0 \quad \hat{\theta}_{ridge}^{noise} = V \, diag\left(\frac{1}{\sigma_1}, \cdots, \frac{1}{\sigma_d}\right) U^t \varepsilon$$

$$\text{high variance}$$

$$\mathbb{E}_{\varepsilon} \| \hat{\theta}_{ridge}^{noise} \|^2 = \sum_{j=1}^{d} \left(\frac{\sigma_j}{\sigma_j^2 + \lambda}\right)^2 \sigma^2$$

$$\text{Variance decreases with } \lambda.$$

Observe: $\quad \lambda$ trades off between bias & variance

Justification of ridge regression estimate $\hat{\theta}_{ridge}$:

Let $X \in \mathbb{R}^{n \times d}$, $y \in \mathbb{R}^n$.

By formula above

$$\hat{\theta}_{ridge} = \left(X^t X + \lambda I_d\right)^{-1} X^t y = \left(X^t X + \lambda I_d\right)^{-1} X^t (X\theta^* + \varepsilon)$$

Let $X = U \Sigma V^t$ be the SVD of $X$, where

$\quad U$ — $n \times d$ matrix with orthonormal columns
$\quad V$ — $d \times d$ matrix with orthonormal columns
$\quad \Sigma$ — $d \times d$ diagonal matrix $= diag(\sigma_1, \cdots, \sigma_d)$ w/ $\sigma_i \geqslant \sigma_{i+1} \geqslant 0$

Note $\quad X^t X = V \Sigma^t \underbrace{U^t U}_{U^t U = I_d} \Sigma V^t = V \underbrace{\Sigma^t I_d \Sigma}_{\Sigma^2} V^t = V \Sigma^2 V^t$

So
$$\hat{\theta}_{ridge} = \left(V\Sigma^2 V^t + \lambda I\right)^{-1}\left[X^t X\, \theta^* + X^t \varepsilon\right]$$

$$= \left(V\Sigma^2 V^t + \lambda I\right)^{-1}\left[V\Sigma^2 V^t \theta^* + V\Sigma^t U^t \varepsilon\right]$$

$I = VV^t$
$$= \left(V(\Sigma^2 + \lambda I)V^t\right)^{-1}\left[V\Sigma^2 V^t \theta^* + V\Sigma^t U^t \varepsilon\right]$$

$$= V(\Sigma^2 + \lambda I)^{-1}V^t\left[V\Sigma^2 V^t \theta^* + V\Sigma^t U^t \varepsilon\right]$$

$$= V(\Sigma^2 + \lambda I)^{-1}\left[\Sigma^2 V^t\theta^* + \Sigma U^t \varepsilon\right]$$

$$= V(\Sigma^2 + \lambda I)^{-1}\Sigma^2 V^t \theta^*$$
$$+ V(\Sigma^2 + \lambda I)^{-1}\Sigma U^t \varepsilon$$

Note $(\Sigma^2 + \lambda I)^{-1} = diag\left(\frac{1}{\sigma_1^2 + \lambda}, \cdots, \frac{1}{\sigma_d^2 + \lambda}\right)$

So
$$\boxed{\begin{aligned}\hat{\theta}_{ridge} &= V\, diag\left(\frac{\sigma_1^2}{\sigma_1^2 + \lambda}, \cdots, \frac{\sigma_d^2}{\sigma_d^2 + \lambda}\right)V^t \theta^* \\ &+ V\, diag\left(\frac{\sigma_1}{\sigma_1^2 + \lambda}, \cdots, \frac{\sigma_d}{\sigma_d^2 + \lambda}\right)U^t \varepsilon\end{aligned}}$$