# COMS 4721: Machine Learning for Data Science
## Lecture 3, 1/24/2017

Prof. John Paisley

Department of Electrical Engineering
& Data Science Institute

Columbia University

# REGRESSION: PROBLEM DEFINITION

### Data

*(extend by 1 dimension for the offset)*

*(real-valued o/p — what makes it a regression problem.)*

Measured pairs $(x, y)$, where $x \in \mathbb{R}^{d+1}$ (input) and $y \in \mathbb{R}$ (output)

### Goal

Find a function $f : \mathbb{R}^{d+1} \to \mathbb{R}$ such that $y \approx f(x; w)$ for the data pair $(x, y)$.
$f(x; w)$ is the *regression function* and the vector $w$ are its parameters.

### Definition of linear regression

A regression method is called *linear* if the prediction $f$ is a linear function of
the unknown parameters $w$.

*(doesn't actually mean we have a linear function of the i/ps x.)*

# Least Squares (continued)

# LEAST SQUARES LINEAR REGRESSION

$$X \to \begin{array}{l} x_i \text{ along the rows} \\ (n, d+1) \end{array}$$

## Least squares solution

Least squares finds the $w$ that minimizes the sum of squared errors. The least squares objective in the most basic form where $f(x; w) = x^T w$ is

$$\mathcal{L} = \sum_{i=1}^{n} (y_i - x_i^T w)^2 \;=\; \|y - Xw\|^2 \;=\; (y - Xw)^T(y - Xw).$$

*matrix*

*dot product btw i/p vector & coefficient $w$*

*vector of errors*

*dot product of error vector with itself we are summing the squares of those errors.*

We defined $y = [y_1, \ldots, y_n]^T$ and $X = [x_1, \ldots, x_n]^T$.

Taking the gradient with respect to $w$ and setting to zero, we find that

$$\nabla_w \mathcal{L} = 2X^T X w - 2X^T y = 0 \quad \Rightarrow \quad w_{\text{LS}} = (X^T X)^{-1} X^T y.$$

In other words, $w_{\text{LS}}$ is the vector that minimizes $\mathcal{L}$.

*(find vector $w$ at which gradient of this function is equal to 0.)*

# PROBABILISTIC VIEW

.

- ► Last class, we discussed the geometric interpretation of least squares.

- ► Least squares also has an insightful probabilistic interpretation that allows us to analyze its properties.

- ► That is, given that we pick this model as reasonable for our problem, we can ask: What kinds of assumptions are we making?

# PROBABILISTIC VIEW

Assuming that the covariance matrix is diagonal, we write covariance as sigma(variance) squared times an identity matrix.

### Recall: Gaussian density in *n* dimensions (σ - variance)

Assume a diagonal covariance matrix $\Sigma = \sigma^2 I$. The density is

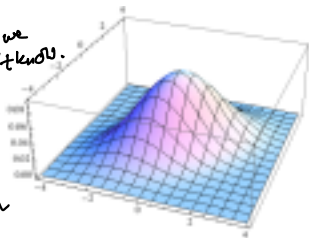$$p(y|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} \exp\left(-\frac{1}{2\sigma^2}(y - \mu)^T (y - \mu)\right). \quad \textcircled{1}$$

output

formed by taking our I/PS and putting them along the rows of the matrix.

$y$ is random variable

What if we restrict the mean to $\mu = Xw$ which we don't know. and find the *maximum likelihood* solution for $w$?

$X$ we know
$w$ we don't know & we want to learn it.

$\mu$ is restricted to have this form

So we plug this matrix vector product $(xw)$ into the Gaussian ① for $\mu$. And then we say let's find the max. likelihood sol^n to $w$.

*Find the value of w that maximises the log-likelihood of the o/p vector, given a mean of this form, diagonal covariance, where each dimension has variance $\sigma^2$.*

## Maximum likelihood for Gaussian linear regression

Plug $\mu = Xw$ into the multivariate Gaussian distribution and solve for $w$ using maximum likelihood.

$$w_{\text{ML}} = \arg\max_w \ln p(y|\mu = Xw, \sigma^2)$$

*And now our goal is to find the vector w that maximises this func. which is a vector that maximises the likelyhood of the data o/p s Y that we see given the corresponding i/ps X.*

$$\arg\max_w -\frac{1}{2\sigma^2}\|y - Xw\|^2 - \frac{n}{2}\ln(2\pi\sigma^2).$$

*because This term doesn't involve w we don't even bother to write it*

<u>Least squares (LS)</u> and <u>maximum likelihood (ML)</u> share the same solution:

$$\text{LS: } \arg\min_w \|y - Xw\|^2 \quad \Leftrightarrow \quad \text{ML: } \arg\max_w -\frac{1}{2\sigma^2}\|y - Xw\|^2$$

*Intuition: The least-square solution corresponds to the maximum likelihood sol$^n$ of this multi-variate Gaussian assumption on our data.*

# PROBABILISTIC VIEW

So therefore in a sense what we can think of L.S. as doing is making an independent Gaussian noise assumption on the error. So this is some intuition that we can develop where we say that the L.S. sol^n corresponds to the max. likelihood sol^n of multivariate Gaussian assumption on our data.

▶ Therefore, in a sense we are making an *independent Gaussian noise* assumption about the error, $\epsilon_i = y_i - x_i^T w$.

▶ Other ways of saying this:

→ Each of the errors are iid as 0 mean Gaussian with some $\sigma^2$. We are effectively modelling our o/p as being equal to the dot product of the i/p with weight w + iid Gaussian noise.

    1) $y_i = x_i^T w + \epsilon_i$,   $\epsilon_i \overset{iid}{\sim} N(0, \sigma^2)$,   for $i = 1, \ldots, n$,

    2) $y_i \overset{ind}{\sim} N(x_i^T w, \sigma^2)$,   for $i = 1, \ldots, n$, ✳

    3) $y \sim N(Xw, \sigma^2 I)$, as on the previous slides. ✳✳

▶ Can we use this probabilistic line of analysis to better understand the maximum likelihood (i.e., least squares) solution?

This like an assumption we're going to make. We're going to say that our noise is i.i.d Gaussian, and that max. likelihood for this modelling assumption is the least squares solution.

* Another way to say, we're modelling our o/ps as independent Gaussian random variables where the mean is equal to the dot product btw the i/p $x_i$ and coefficient vector w.

** As the previous slide we can simply say that we've making the assumption that vector Y is multivaried Gaussian with mean $= Xw$ and covariance $\sigma^2 I$

         i/p along each row

We're going to say that our noise is IID Gaussian. Essentially and that max. likelihood for the modelling assumption is the least-squared sol$^n$.

# PROBABILISTIC VIEW

### Expected solution

*we observe are generated*
*Gaussian with mean*
*$Xw$*

*S/ps thet* [↑]

**Given:** The *modeling assumption* that $y \sim N(Xw, \sigma^2 I)$. *iid gaussian noise*

We can calculate the <u>expectation of the ML solution under this distribution</u>,

*nothing random there*

$$\mathbb{E}[w_{ML}] = \underline{\mathbb{E}[(X^T X)^{-1} X^T y]} \quad \left( = \int \left[ (X^T X)^{-1} X^T y \right] p(y|X,w)\,dy \right)^{*}$$

*we're making an distribution*
*assumption on how y is*
*generated - so y is random.*

$$= (X^T X)^{-1} X^T \mathbb{E}[y]$$
$$= (X^T X)^{-1} X^T X w$$
$$= w$$

Therefore $w_{ML}$ is an *unbiased* estimate of $w$, i.e., $\mathbb{E}[w_{ML}] = w$.

*we don't know w. We are*
*assuming, there exists some*
*w.*

This is the expectation function. So, of course every expectation you have some assumption of the distribution, you're taking that expectation over.

In this case, it's distribution on y. So you have a function of y times the distribution of y, and now we're integrating over y.

So what this is saying intuitively is that if we have some ground truth value for W. And we generate, we have some input, inputs X that we construct in a matrix in this way.And then we generate an output Y according to this distribution.
And then using that output y, we solve the maximum likelihood solution for w. So we have the true w, and then using the random vector y, solve the maximum likelihood solution for w, the expectation of our maximum likelihood solution is equal to the truth.
This case what were saying is the maximum likelihood solution for the vector w is an unbiased estimate.
Of the ground truth vector w, which we don't have access to.
So this is good, least squares or maximum likelihood for this model is going to in expectation, give us the true parameter, which is what we're trying to learn. I can't really understand the big thing,

we just added a random variable with expectation equal to zero. And the expectation didn't change.

▶ Even though the "expected" maximum likelihood solution is the correct one, should we actually expect to get something near it?

If I have a Gaussian random variable with mean $\mu$ & variance $\sigma^2$. And $\sigma^2$ is huge, then even though that random variable in expectation is equal to the $\mu$. Since the variance is so huge, I don't actually expect to see something close to the mean.

So, we have shown that maximum likelihood sol$^n$ is the correct one. But now we have calculate the variance of that sol$^n$ under the same modelling assumption.

▶ Even though the "expected" maximum likelihood solution is the correct one, should we actually expect to get something near it?

▶ We should also look at the covariance. Recall that if $y \sim N(\mu, \Sigma)$, then

$$\text{Var}[y] = \mathbb{E}[(y - \mathbb{E}[y])(y - \mathbb{E}[y])^T] = \Sigma.$$

covariance = expectation of
the outer product
with it's subtracted
off.

$$\left[ \begin{array}{c} \text{For Gaussian it is} \\ \Sigma. \end{array} \right]$$

# REVIEW: AN EQUALITY FROM PROBABILITY

- Even though the "expected" maximum likelihood solution is the correct one, should we actually expect to get something near it?

- We should also look at the covariance. Recall that if $y \sim N(\mu, \Sigma)$, then

$$\text{Var}[y] = \mathbb{E}[(y - \mathbb{E}[y])(y - \mathbb{E}[y])^T] = \Sigma.$$

- Plugging in $\mathbb{E}[y] = \mu$ (expectation), this is equivalently written as

$$
\begin{aligned}
\text{Var}[y] &= \mathbb{E}[(y - \mu)(y - \mu)^T] \quad \text{take the outer product} \\
&= \mathbb{E}[yy^T - y\mu^T - \mu y^T + \mu\mu^T] \\
&= \mathbb{E}[yy^T] - \mu\mu^T \quad \mathbb{E}(y) = \mu
\end{aligned}
$$

- Immediately, we also get $\mathbb{E}[yy^T] = \Sigma + \mu\mu^T$.

with this gaussian

# PROBABILISTIC VIEW

*Return to LS linear regression problem*

*And*

*Calculate the covari ance of the max.*

### Variance of the solution

*likelihood sol² under the Gaussian*

Returning to least squares linear regression, we wish to find

*max. likelihood sol² — It's expectation*

$$\text{Var}[w_{\text{ML}}] = \mathbb{E}[(w_{\text{ML}} - \mathbb{E}[w_{\text{ML}}])(w_{\text{ML}} - \mathbb{E}[w_{\text{ML}}])^T]$$

$$= \mathbb{E}[w_{\text{ML}} w_{\text{ML}}^T] - \mathbb{E}[w_{\text{ML}}]\mathbb{E}[w_{\text{ML}}]^T.$$

*⤷ From the previous slide where the variance of random vector is equal to the expectation of the outer product subtracted by the outer product of its mean. (○)*

---

[1] Aside: For matrices $A$, $B$ and vector $c$, recall that $(ABc)^T = c^T B^T A^T$.

.

### Variance of the solution

Returning to least squares linear regression, we wish to find

$$
\begin{aligned}
\mathrm{Var}[w_{\mathrm{ML}}] &= \mathbb{E}[(w_{\mathrm{ML}} - \mathbb{E}[w_{\mathrm{ML}}])(w_{\mathrm{ML}} - \mathbb{E}[w_{\mathrm{ML}}])^T] \\
&= \mathbb{E}[w_{\mathrm{ML}} w_{\mathrm{ML}}^T] - \mathbb{E}[w_{\mathrm{ML}}]\mathbb{E}[w_{\mathrm{ML}}]^T.
\end{aligned}
$$

The sequence of equalities follows:[1]

$$
\check{\mathrm{Var}}[w_{\mathrm{ML}}] = \mathbb{E}[(X^TX)^{-1}X^Tyy^TX(X^TX)^{-1}] - ww^T
$$

*(handwritten annotations:)*
→ expectation of the max. likelihood sol$^n$ is equal to the truth (already shown).

max. likelihood sol$^n$ for w.

[whatever w is the expectation of ML is equal to w]

even though we don't know what it is. We can simply plug that value in there as a placeholder.

---

[1] Aside: For matrices $A$, $B$ and vector $c$, recall that $(ABc)^T = c^TB^TA^T$.

# PROBABILISTIC VIEW

## Variance of the solution

Returning to least squares linear regression, we wish to find

$$
\begin{aligned}
\mathrm{Var}[w_{\mathrm{ML}}] &= \mathbb{E}[(w_{\mathrm{ML}} - \mathbb{E}[w_{\mathrm{ML}}])(w_{\mathrm{ML}} - \mathbb{E}[w_{\mathrm{ML}}])^T] \\
&= \mathbb{E}[w_{\mathrm{ML}} w_{\mathrm{ML}}^T] - \mathbb{E}[w_{\mathrm{ML}}]\mathbb{E}[w_{\mathrm{ML}}]^T.
\end{aligned}
$$

The sequence of equalities follows:[1]

$$
\begin{aligned}
\mathrm{Var}[w_{\mathrm{ML}}] &= \mathbb{E}[(X^T X)^{-1} X^T y y^T X (X^T X)^{-1}] - w w^T \\
&= \underbrace{(X^T X)^{-1} X^T} \mathbb{E}[y y^T] X (X^T X)^{-1} - w w^T
\end{aligned}
$$

*matrices involving $X$ are contant we can bring them out of the expectation.*

---

[1] Aside: For matrices *A*, *B* and vector *c*, recall that $(ABc)^T = c^T B^T A^T$.

# PROBABILISTIC VIEW

### Variance of the solution

Returning to least squares linear regression, we wish to find

$$
\begin{aligned}
\text{Var}[w_{\text{ML}}] &= \mathbb{E}[(w_{\text{ML}} - \mathbb{E}[w_{\text{ML}}])(w_{\text{ML}} - \mathbb{E}[w_{\text{ML}}])^T] \\
&= \mathbb{E}[w_{\text{ML}}w_{\text{ML}}^T] - \mathbb{E}[w_{\text{ML}}]\mathbb{E}[w_{\text{ML}}]^T.
\end{aligned}
$$

The sequence of equalities follows:[1]

$$
\begin{aligned}
\text{Var}[w_{\text{ML}}] &= \mathbb{E}[(X^TX)^{-1}X^Tyy^TX(X^TX)^{-1}] - ww^T \\
&= (X^TX)^{-1}X^T\mathbb{E}[yy^T]X(X^TX)^{-1} - ww^T \\
&= (X^TX)^{-1}X^T(\sigma^2I + Xww^TX^T)X(X^TX)^{-1} - ww^T
\end{aligned}
$$

*[handwritten annotation: we assume w has this relationship with o/p so can plug this function in here]*

*[handwritten annotation: is eqal to the covariance of y + outer product of mean of vector y which is eqal to Xw.]*

---

[1] Aside: For matrices $A$, $B$ and vector $c$, recall that $(ABc)^T = c^TB^TA^T$.

# PROBABILISTIC VIEW

### Variance of the solution

Returning to least squares linear regression, we wish to find

$$
\begin{aligned}
\text{Var}[w_{\text{ML}}] &= \mathbb{E}[(w_{\text{ML}} - \mathbb{E}[w_{\text{ML}}])(w_{\text{ML}} - \mathbb{E}[w_{\text{ML}}])^T] \\
&= \mathbb{E}[w_{\text{ML}} w_{\text{ML}}^T] - \mathbb{E}[w_{\text{ML}}]\mathbb{E}[w_{\text{ML}}]^T.
\end{aligned}
$$

The sequence of equalities follows:[1]

$$
\begin{aligned}
\text{Var}[w_{\text{ML}}] &= \mathbb{E}[(X^TX)^{-1}X^T yy^T X(X^TX)^{-1}] - ww^T \\
&= (X^TX)^{-1}X^T \mathbb{E}[yy^T] X(X^TX)^{-1} - ww^T \\
&= \underbrace{(X^TX)^{-1}X^T}(\underline{\sigma^2 I} + Xww^TX^T)\underbrace{X(X^TX)^{-1}} - ww^T
\end{aligned}
$$

separate out = $(X^TX)^{-1}X^T \sigma^2 I X(X^TX)^{-1} + \cdots$

the matrices $(X^TX)^{-1}X^T Xww^TX^TX(X^TX)^{-1} - ww^T$

---

[1] Aside: For matrices $A$, $B$ and vector $c$, recall that $(ABc)^T = c^T B^T A^T$.

### Variance of the solution

Returning to least squares linear regression, we wish to find

$$
\begin{aligned}
\text{Var}[w_{\text{ML}}] &= \mathbb{E}[(w_{\text{ML}} - \mathbb{E}[w_{\text{ML}}])(w_{\text{ML}} - \mathbb{E}[w_{\text{ML}}])^T] \\
&= \mathbb{E}[w_{\text{ML}} w_{\text{ML}}^T] - \mathbb{E}[w_{\text{ML}}]\mathbb{E}[w_{\text{ML}}]^T.
\end{aligned}
$$

The sequence of equalities follows:[1]

$$
\begin{aligned}
\text{Var}[w_{\text{ML}}] &= \mathbb{E}[(X^TX)^{-1}X^T yy^T X(X^TX)^{-1}] - ww^T \\
&= (X^TX)^{-1}X^T \mathbb{E}[yy^T] X(X^TX)^{-1} - ww^T \\
&= (X^TX)^{-1}X^T (\sigma^2 I + Xww^TX^T) X(X^TX)^{-1} - ww^T \\
&= (X^TX)^{-1}X^T \sigma^2 I X(X^TX)^{-1} + \cdots \\
&= \underbrace{(X^TX)^{-1}X^T X}ww^T \underbrace{X^T X(X^TX)^{-1}} - ww^T \\
&= \sigma^2(X^TX)^{-1}
\end{aligned}
$$

$\sigma^2 (x^T x)^{-1} x^T x (x^T x)^{-1}$

$= \sigma^2 (x^T x^{-1})$

$\mathcal{I}$

$I \searrow ww^T - ww^T = 0$

---

# PROBABILISTIC VIEW

*If we make a Gaussian assumption, we assume that the vector of responses y is equal to matrix of feature vectors X times weight vector w + independent Gaussian noise $\sigma^2$*

*Note: we have only defined a distribution of y, we haven't defined any distribution on x.*

▶ We've shown that, under the Gaussian assumption $y \sim N(Xw, \sigma^2 I)$,

*And we solve LS or max. likelihood*
*Sol^n for w, we assume*

$$\mathbb{E}[w_{ML}] = w, \quad \text{Var}[w_{ML}] = \sigma^2(X^T X)^{-1}.$$

*there some true w. But we observe*
*y which is a random vector acc. to $N(Xw, \sigma \mathbb{I}^{-1})$ distribution. We try to do the*
*inverse of learning what is w? Using max. likelihood.*

▶ When there are very large values in $\sigma^2(X^T X)^{-1}$, the values of $w_{ML}$ are very sensitive to the measured data $y$ (more analysis later).

*We don't trust max. likelihood sol^n when the values in covariance are large.*

▶ This is bad if we want to analyze and predict using $w_{ML}$.

*↳ Then the expected value of our max. likelihood sol^n is equal to the true.*
*But covariance is $\sigma^2(X^T X)^{-1}$.*

*when the values in the covariance are v.v. large we can't say the max. likelihood sol^n is close to the truth.*

# Ridge Regression

*Why does penalizing the magnitude of w reduce the variance?*

▶ We saw how with least squares, the values in $w_{\text{ML}}$ may be huge.

*and deviate wildly from the ground truth depending on what the matrix $(XX^T)^{-1}$ looks like.*

▶ In general, when developing a model for data we often wish to *constrain* the model parameters in some way.

*In this case, it might take the form of penalizing values of w that we might consider too large.*

▶ There are many models of the form

*penalty term*

$$w_{\text{OPT}} = \arg\min_w \|y - Xw\|^2 + \lambda g(w).$$

*function of our model parameters*

*w*

*sum of squared errors of our approximation*

▶ The added terms are

1. $\lambda > 0$ : a regularization parameter,
2. $g(w) > 0$ : a penalty function that encourages desired properties about $w$.

*that penalizes values of the vector w in some way, in order to encourage properties of w that we might want in advance.*

# RIDGE REGRESSION

*naming a model of a particular form, where we use a particular regularization function g.*

Ridge regression is one $g(w)$ that addresses variance issues with $w_{ML}$.

It uses the squared penalty on the regression coefficient vector $w$,

$$w_{RR} = \arg\min_w \|y - Xw\|^2 + \lambda\|w\|^2$$

The term $g(w) = \|w\|^2$ penalizes large values in $w$.

*g is squared magnitude of the vector w.*

However, there is a *tradeoff* between the first and second terms that is controlled by $\lambda$. *(parameter)*

▶ Case $\lambda \to 0$ : $w_{RR} \to w_{LS}$ *($II^{nd}$ term disappears).*

▶ Case $\lambda \to \infty$ : $w_{RR} \to \vec{0}$ *(any non-negative value of w has essentially infinite penalty.)*

*vector of all zeros.*

# RIDGE REGRESSION SOLUTION

**Objective:** We can solve the ridge regression problem using exactly the same procedure as for least squares,

$$\mathcal{L} = \|y - Xw\|^2 + \lambda\|w\|^2$$
$$= (y - Xw)^T(y - Xw) + \lambda w^T w.$$

**Solution:** First, take the gradient of $\mathcal{L}$ with respect to $w$ and set to zero,

$$\nabla_w \mathcal{L} = -2X^T y + 2X^T Xw + 2\lambda w = 0 \quad \left(\begin{array}{l}\text{In this case, we can't}\\ \text{solve in closed form.}\end{array}\right)$$

Then, solve for $w$ to find that

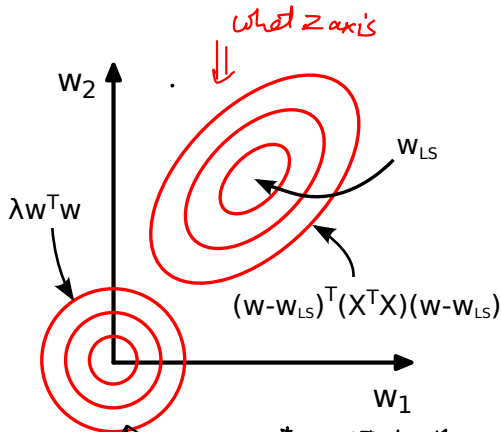$$w_{\text{RR}} = (\lambda I + X^T X)^{-1} X^T y.$$

$\lambda \geq 0$    least squares sol^n

$\lambda \to \infty$, the inverse of that matrix goi. to be a matrix of all zeros

# RIDGE REGRESSION GEOMETRY

There is a tradeoff between squared error and penalty on *w*.

We can write both in terms of *level sets*: Curves where function evaluation gives the same number.

The sum of these gives a new set of levels with a unique minimum.

*what z axis*

$w_2$

$w_{LS}$

$\lambda w^T w$

$(w - w_{LS})^T (X^T X)(w - w_{LS})$

$w_1$

*same value at any pt along the circle*

You can check that we can write: *involves w*

*does it? nah!*

$$\|y - Xw\|^2 + \lambda\|w\|^2 = (w - w_{LS})^T(X^TX)(w - w_{LS}) + \lambda w^T w + (\text{const. w.r.t. } w)$$

*write these 2 functions as their own independent level sets.*

*Means if you evaluate a function at any pt. along the curve, the sol^n has the*

# DATA PREPROCESSING

*penalizes each dimension of w equally.* [λ||w||²]

[0-1, 1-100 diff. dimensions.
Didn't have this problem with
least squares, as we had just
no restriction on w.]

Ridge regression is one possible regularization scheme. For this problem, we first assume the following *preprocessing* steps are done:

[If two dimensions are
equally important, Ridge
regression is going to
penalize on the dimension
that takes greater
values.]

1. The mean is subtracted off of $y$:

$$y \leftarrow y - \frac{1}{n} \sum_{i=1}^{n} y_i.$$

2. The dimensions of $x_i$ have been *standardized* before constructing $X$:

$$x_{ij} \leftarrow (x_{ij} - \bar{x}_{\cdot j})/\hat{\sigma}_j, \quad \hat{\sigma}_j = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (x_{ij} - \bar{x}_{\cdot j})^2}.$$

i.e., subtract the empirical mean and divide by the empirical standard deviation for each dimension.

3. We can show that there is no need for the dimension of 1's in this case.

↓ why is it no needed?

# Some Analysis of Ridge Regression

## RIDGE REGRESSION VS LEAST SQUARES

The solutions to least squares and ridge regression are clearly very similar,

$$w_{LS} = (X^T X)^{-1} X^T y \quad \Leftrightarrow \quad w_{RR} = (\lambda I + X^T X)^{-1} X^T y.$$

Discuss what this $\lambda$ is doing.

▶ We can use linear algebra and probability to compare the two.

▶ This requires the *singular value decomposition*, which we review next.

# REVIEW: SINGULAR VALUE DECOMPOSITIONS

*we have more observations than dimensions in our problem*

▶ We can write any $n \times d$ matrix $X$ (assume $n > d$) as $X = USV^T$, where
   *each column is unit length & orthogonal to every other column*
   1. $U$: $n \times d$ and orthonormal in the columns, i.e. $U^T U = I$.
   2. $S$: $d \times d$ non-negative diagonal matrix, i.e. $S_{ii} \geq 0$ and $S_{ij} = 0$ for $i \neq j$.
   3. $V$: $d \times d$ and orthonormal, i.e. $\underbrace{V^T V = VV^T}_{\text{since it is square}} = I$. $\left[\begin{matrix} \text{Both columns & rows of } V \text{ are} \\ \text{orthonormal.} \end{matrix}\right]$
   *square*

▶ From this we have the immediate equalities
   *replace byte side for X.*
   $$X^T X = (USV^T)^T (USV^T) = VS^2 V^T, \quad XX^T = US^2 U^T.$$
   $$VSU^T USV^T$$

▶ Assuming $S_{ii} \neq 0$ for all $i$ (i.e., "$X$ is full rank"), we also have that

   $$(X^T X)^{-1} = (VS^2 V^T)^{-1} = VS^{-2} V^T.$$ *that's how can the inverse is*

   Proof: Plug in and see that it satisfies definition of inverse

   $$(X^T X)(X^T X)^{-1} = VS^2 V^T VS^{-2} V^T = I.$$

   $$VS^{-2} V^T$$

# LEAST SQUARES AND THE SVD

*Using SVD to analyze the least squares solution.*

Using the SVD we can rewrite the variance, *replace with svd*

$$\mathrm{Var}[w_{\mathrm{LS}}] = \sigma^2(X^TX)^{-1} = \sigma^2 V S^{-2} V^T.$$

*→ square, become close to zero. then invert very large values.*

This <u>inverse</u> becomes <u>huge</u> when $S_{ii}$ is very small for some values of $i$.
(Aside: This happens when columns of $X$ are highly correlated.)

*hence singular values that are very small.*

The least squares prediction for new data is

$$y_{\mathrm{new}} = x_{\mathrm{new}}^T w_{\mathrm{LS}} = x_{\mathrm{new}}^T (X^TX)^{-1}X^T y = x_{\mathrm{new}}^T \underbrace{V S^{-1} U^T}_{\text{svd representation}} y.$$

*→ very huge, for small values*

When $S^{-1}$ has very large values, this can lead to unstable predictions.

*depending upon how the vector x correlates with these singular vectors that small singular values.*

# RIDGE REGRESSION VS LEAST SQUARES I

Manipulating Ridge Regression sol$^n$ in a certain way to relate it to the least squares sol$^n$.

## Relationship to least squares solution

Recall for two symmetric matrices, $(AB)^{-1} = B^{-1}A^{-1}$.

$$
\begin{aligned}
w_{\text{RR}} &= (\lambda I + X^TX)^{-1}X^Ty \qquad \overset{n^{\frac{3}{2}}}{\underset{}{}} \quad \text{(multiply \& divide by} \\
&\qquad\qquad\qquad\qquad\qquad\qquad\qquad \text{the same thing.)} \\
&= (\lambda I + X^TX)^{-1}(X^TX)\underbrace{(X^TX)^{-1}X^Ty}_{w_{\text{LS}}} \\
&= [(X^TX)(\lambda(X^TX)^{-1} + I)]^{-1}(X^TX)w_{\text{LS}} \\
&= (\lambda(X^TX)^{-1} + I)^{-1}(X^TX)^{-1}(X^TX)w_{\text{LS}} \\
&= (\lambda(X^TX)^{-1} + I)^{-1}w_{\text{LS}} \qquad \left[\begin{array}{l}\text{we can simply manipulate the} \\ \text{least square Sol$^n$ acc. to this} \\ \text{vector to}\end{array}\right.
\end{aligned}
$$

Can use this to prove that the solution shrinks toward zero: $\|w_{\text{RR}}\|_2 \leq \|w_{\text{LS}}\|_2$. get

We should be able to expect the ridge regression sol$^n$ to have smaller magnitude then the least squares solution.

Adding something that's non-negative singular values to it. So you can imagine ridge regression sol$^n$? this as dividing by something that is $>1$. So you're shrinking this.

# RIDGE REGRESSION VS LEAST SQUARES II

*Replace $x$ by its svd.*

Continue analysis with the SVD: $X = USV^T \rightarrow (X^TX)^{-1} = VS^{-2}V^T$:

$$
\begin{aligned}
w_{RR} &= (\lambda(X^TX)^{-1} + I)^{-1}w_{LS} \\
V \rightarrow \text{square orthogonal matrix} \quad &= (\lambda VS^{-2}V^T + I)^{-1}w_{LS} \\
I, S^{-2} \rightarrow \text{diagonal matrix} \quad &= V(\lambda S^{-2} + I)^{-1}V^Tw_{LS} \\
&:= V\stackrel{\downarrow}{M}V^Tw_{LS}
\end{aligned}
$$

$M$ is a diagonal matrix with $M_{ii} = \frac{S_{ii}^2}{\lambda + S_{ii}^2}$. We can pursue this to show that

*how?*

*I get this*

$$
w_{RR} = VS_\lambda^{-1}U^Ty, \quad S_\lambda^{-1} = \begin{bmatrix} \frac{S_{11}}{\lambda + S_{11}^2} & & 0 \\ & \ddots & \\ 0 & & \frac{S_{dd}}{\lambda + S_{dd}^2} \end{bmatrix}
$$

$VM\underline{V^TVS^{-1}U^T}y \rightarrow VMS^{-1}U^Ty$

*As $S_{dd} \rightarrow 0$*

*① as $\lambda \rightarrow 0$,*
*goes to $\rightarrow \infty$*

*② now $\lambda > 0$*
*$\rightarrow$ goes to $\rightarrow 0$*

Compare with $w_{LS} = VS^{-1}U^Ty$, which is the case where $\lambda = 0$ above.

*Also, use svd on the soln of least squares*

*[ $\lambda$ is essentially killing off these small singular values.*

# RIDGE REGRESSION VS LEAST SQUARES III

*how RR relates to LS as a LS problem.*

Ridge regression can also be seen as a special case of least squares.

Define $\hat{y} \approx \hat{X}w$ in the following way,

$$d \text{ zeros} \left\{ \begin{bmatrix} y \\ 0 \\ \vdots \\ 0 \end{bmatrix} \approx \begin{bmatrix} - & X & - \\ & \sqrt{\lambda} & & 0 \\ & & \ddots & \\ & 0 & & \sqrt{\lambda} \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix} \right.$$

*Attach a diagonal matrix x along the bottom of X that has $\sqrt{\lambda}$ along the diagonal. AND*

If we solved $w_{LS}$ for *this* regression problem, we find $w_{RR}$ of the *original* problem: Calculating $(\hat{y} - \hat{X}w)^T(\hat{y} - \hat{X}w)$ in two parts gives ·

*Solve the least square sol^n for this problem*

$$\begin{aligned} (\hat{y} - \hat{X}w)^T(\hat{y} - \hat{X}w) &= (y - Xw)^T(y - Xw) + (\sqrt{\lambda}w)^T(\sqrt{\lambda}w) \\ &= \|y - Xw\|^2 + \lambda\|w\|^2 \end{aligned}$$

*So RR is almost like a augmented LS problem.*

$\lambda \to \infty$, sol^n is a vector of 0s
$\lambda \to 0$, sol^n $\to$ LS sol^n

trace of sol^n as a function of $\lambda$.
way to understand how $\lambda$ changes our r r sol^n.

plot of weight vector w.
as a function of degrees of freedom.

Degrees of freedom:

$$df(\lambda) = \text{trace}\left[X(X^TX + \lambda I)^{-1}X^T\right]$$

why? $$= \sum_{i=1}^{d} \frac{S_{ii}^2}{\lambda + S_{ii}^2}$$

Why is it equal to the trace of M?

$df(\lambda)=$ $\lambda$ for which

The sol^n fundamentally changes as $\lambda$ goes from 0 to $\infty$.
Foren: wt of age ↑. So ↑ age
response ↑
wt of age ↓. ↑ age ↓
response ↓
we're saying that the age of person in this example
fundamentally changes for diff. $\lambda$s.

8 values of wt. vector.
for this $\lambda$.

This gives a way of
visualizing relationships.

We will discuss methods for
picking $\lambda$ later.

as a function of degrees of
freedom. that is wt. br ridge
regression sol^n for age dimension.

how do you interpret this?
no of df = dim. (# features)

$\lambda = \infty$

dimensionality of
problem.

no. of
covariates

$\lambda = 0$

df = d