

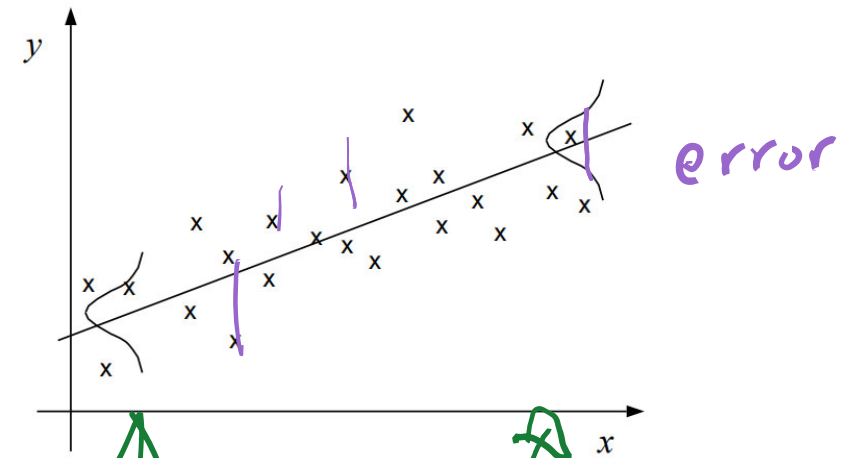
CS 189/289

Today's lecture:

Linear regression (MLE + conditional Gaussians)

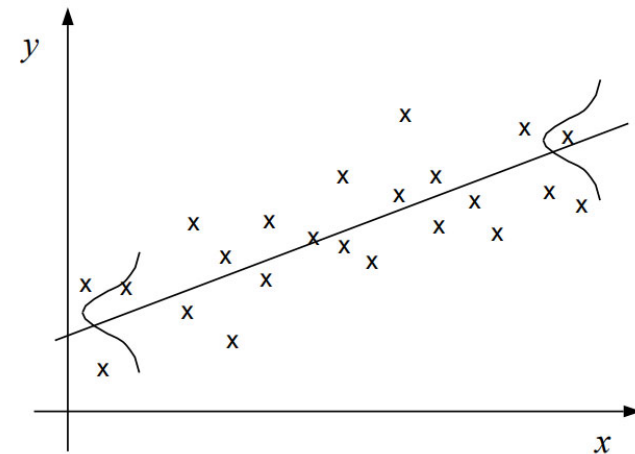
Regression

- *Supervised learning*: data pairs $D = \{(x_i, y_i)\}$, where, x_i may be discrete and/or continuous.
- *Regression*: label, y_i , is a real-valued, e.g., $y_i \in \mathbb{R}$.
- Formally, want $p(y|x)$, the conditional pdf.
- "Point" prediction is then $\hat{y} = E_Y[p(Y|X = x)]$.



Regression examples

- Covid infection rates from zip code and vaccination rate, etc.
- How much a particular protein will bind to a drug target.
- A person's blood pressure from their genetics.
- Tracking - object location in video at the next time-step.
- Housing prices, crime rates, stock prices, *etc.*
- Earliest regression: Legendre in 1805, and Gauss in 1809, both estimating orbits of bodies about the sun.

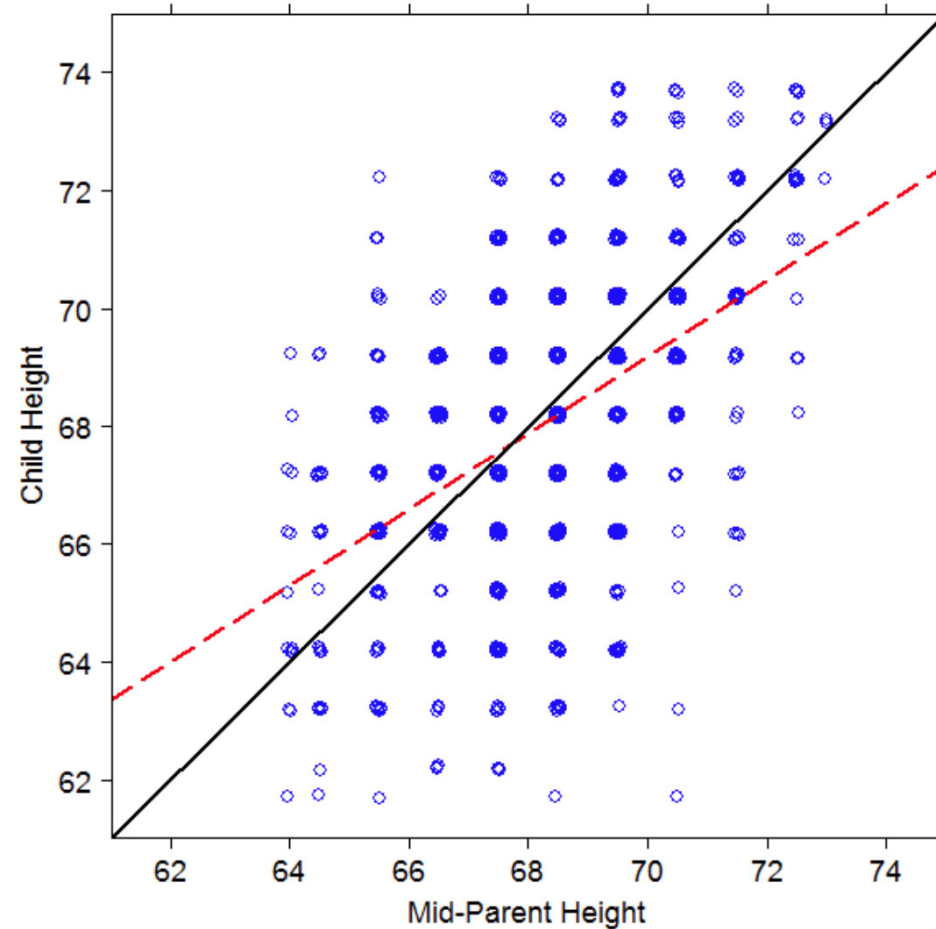


History of the term "Regression"

- Sir Francis Galton (1822-1911) "regression to the mean".

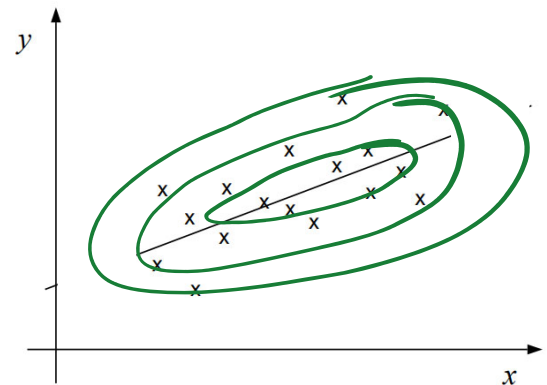


"It appeared from these experiments that the offspring did not tend to resemble their parents in size, but always to be more mediocre than they – to be smaller than the parents, if the parents were large; to be larger than the parents, if the parents were small."



Possible Regression Tactics (to estimate $p(y|x)$)

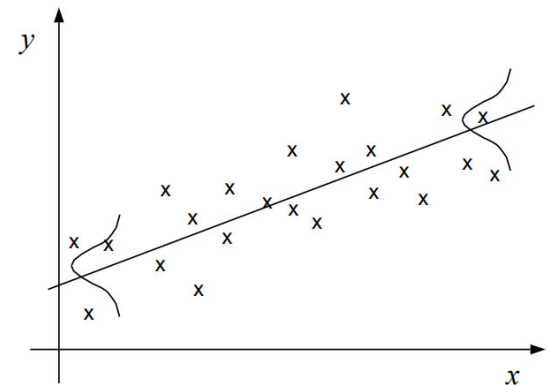
- Our data are drawn from some distribution, $(X, Y) \sim p(x, y)$.
 - What are possible strategies to estimate $p(y|x)$?
1. Estimate $p(x, y|\theta)$ (e.g. MVG for RVs X, Y), and then use fitted model to compute $p(y|x, \hat{\theta}) = \frac{p(y, x|\hat{\theta})}{p(x|\hat{\theta})} = \frac{p(y, x|\hat{\theta})}{\int_y p(y, x|\hat{\theta}) dy}$.



Possible Regression Tactics (to estimate $p(y|x)$)

- Our data are drawn from some distribution, $(X_i, Y_i) \sim p(x, y)$.
 - What are possible strategies to estimate $p(y|x)$?
1. Estimate $p(x, y|\theta)$ (e.g. **MVG for RVs X, Y**), and then use fitted model to compute $p(y|x, \hat{\theta}) = \frac{p(y, x|\hat{\theta})}{p(x|\hat{\theta})} = \frac{p(y, x|\hat{\theta})}{\int_y p(y, x|\hat{\theta}) dy}$.
 2. Consider the inputs to be fixed, and model only the output as a RV. That is, **directly model the conditional $p(y|x, \hat{\theta})$** .

"generative", vs. *"discriminative"*



Linear Regression

- Takes the discriminative approach.
- Predictions are a linear function of the parameters:

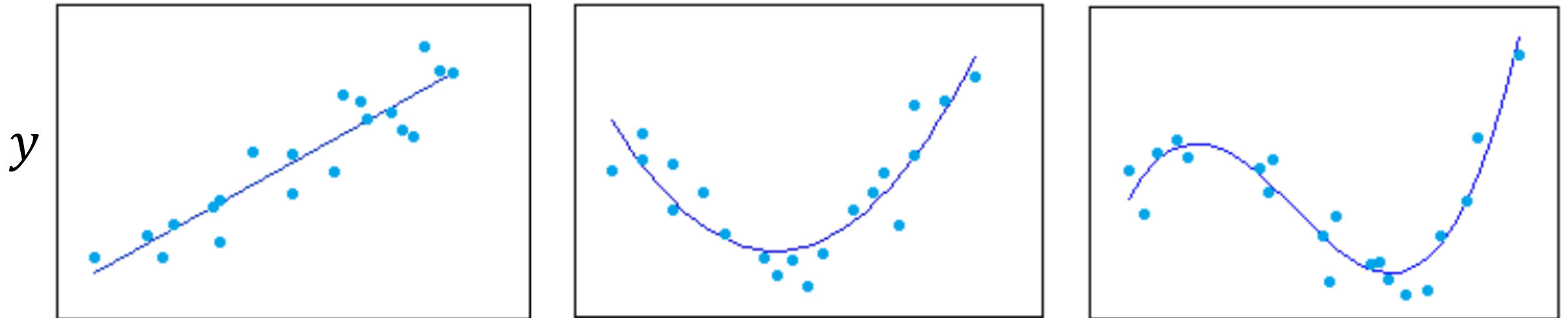
$$\hat{y} = E_Y[p(y|x)] = \mathbf{w}^T \mathbf{x} + \mathbf{w}_0, \text{ for } \mathbf{w}, \mathbf{x} \in \mathbb{R}^d.$$

- \mathbf{w}_0 is called the "offset"/"bias"/"intercept".
- Instead of a bias, we can make an extra feature that is always $\mathbf{1}$.
- Now use $\mathbf{x}' = [\mathbf{x}, \mathbf{1}]$ and $\hat{y} = \mathbf{w}^T \mathbf{x}'$.

How useful can a linear model be?!

Which of these curves can be modelled by linear regression?

$$\hat{y} = E_Y[p(y|x)] = \mathbf{w}^T \mathbf{x} + w_0, \text{ for } \mathbf{w}, \mathbf{x} \in \mathbb{R}^d$$

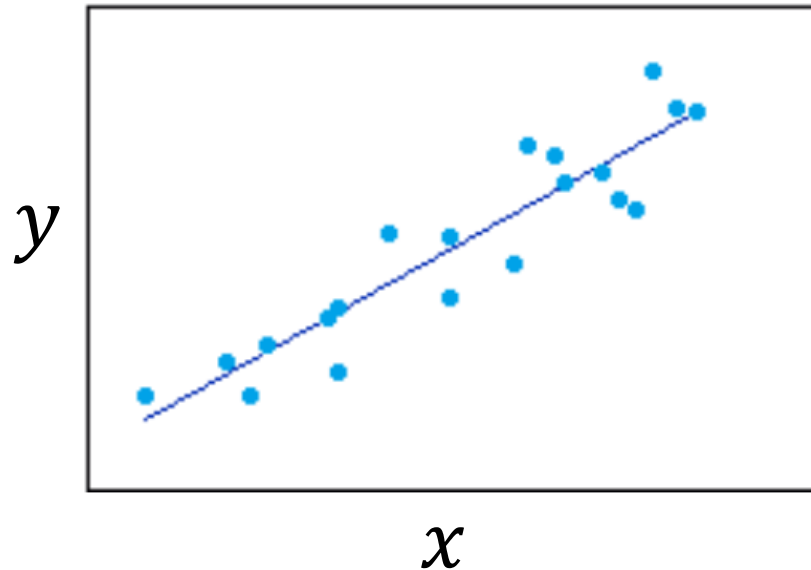


How useful can a linear model be?!

$w, x \in \mathbb{R}^1$

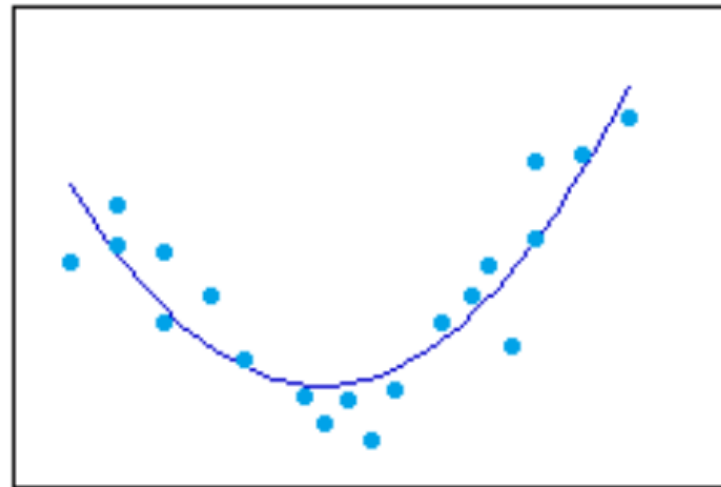
$$\hat{y} = w^T x$$

Linear



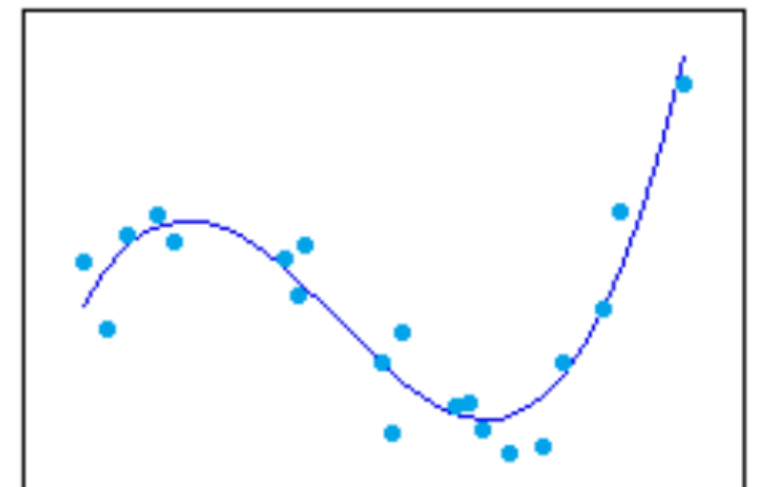
$$\hat{y} = w^T [x, x^2]$$

Quadratic



$$\hat{y} = w^T [x, x^2, x^3]$$

Cubic



For full generality, $x \in \mathbb{R}^D$ need the cross-terms and bias terms, e.g., quadratic $[1, x_1, x_2, x_1^2, x_2^2, x_1 x_2]$.

Basis expansion of raw input space

$$\boldsymbol{x} \in \mathbb{R}^{d=2} = [x_1, x_2] \rightarrow [1, x_1, x_2, x_1 x_2, x_1^2, x_2^2] \in \mathbb{R}^{k=6}$$

Polynomial expansion of order 2 (i.e. quadratic)

- Denote *basis expansion* of the (raw) input features:
 $\Phi(\boldsymbol{x}): \mathbb{R}^d \rightarrow \mathbb{R}^k$.

For $d = 1$, polynomial expansions of order $k = 2, 3$:

- For a quadratic expansion, $\Phi(x) = [1, x, x^2]$, and $k = 2$.
- For a cubic expansion, $\Phi(x) = [1, x, x^2, x^3]$, and $k = 3$.
- For a linear/identity basis expansion, $\Phi(x) = x$, and $k = d$.

Basis expansion of raw input space

Basis functions are pre-determined, so just a notational change:

$$\hat{y} = E_Y[p(y|x)] = \mathbf{w}^T \Phi(\mathbf{x}), \text{ for } \mathbf{w} \in \mathbb{R}^k, \mathbf{x} \in \mathbb{R}^d$$

In this lecture, for simplicity of notation, we will assume that this expansion has already been done, and just write $\hat{y} = \mathbf{w}^T \mathbf{x}$.

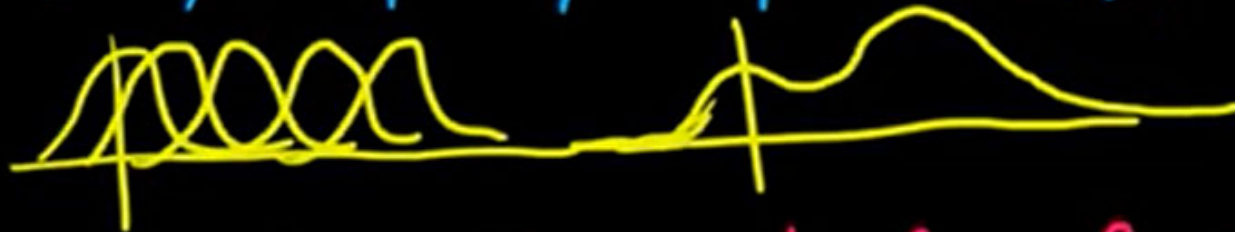
Many basis possible functions!

Polynomials: $f(x) = w_1 + w_2 x(1) + w_3 x(2) + w_4 x(1)^2 + w_5 x(2)^2 + w_6 x(1)x(2)$

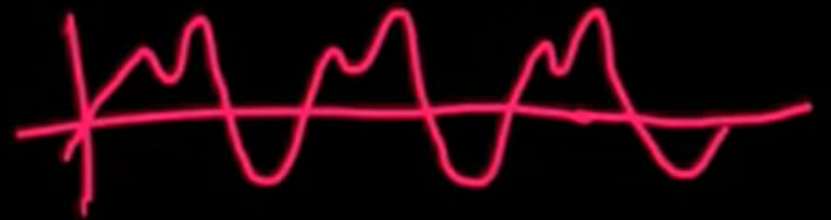
$$\phi(x) = (1, x(1), x(2), x(1)^2, x(2)^2, x(1)x(2))$$

$$\exp\left(\frac{-(x-c)^2}{r^2}\right)$$

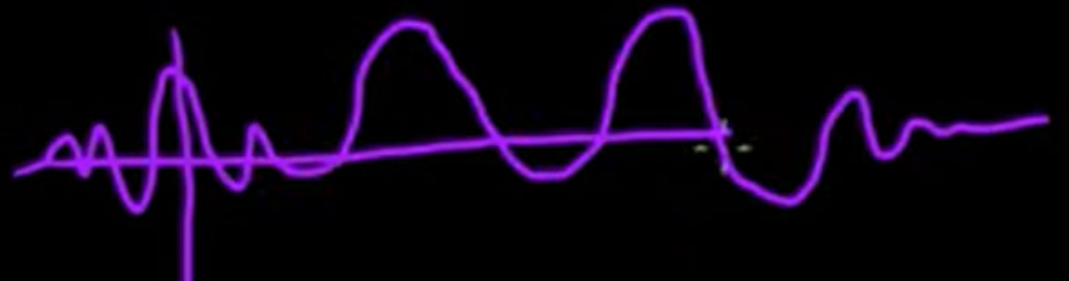
Radial Basis fns:



Fourier basis:

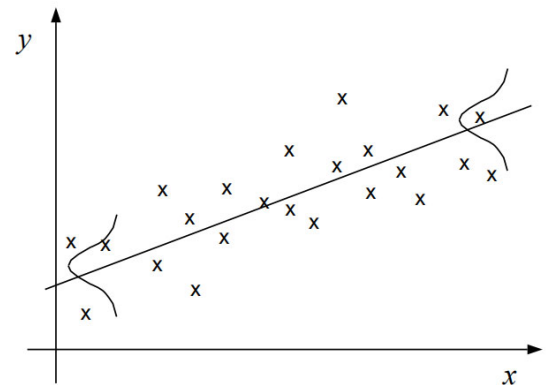


Wavelets:



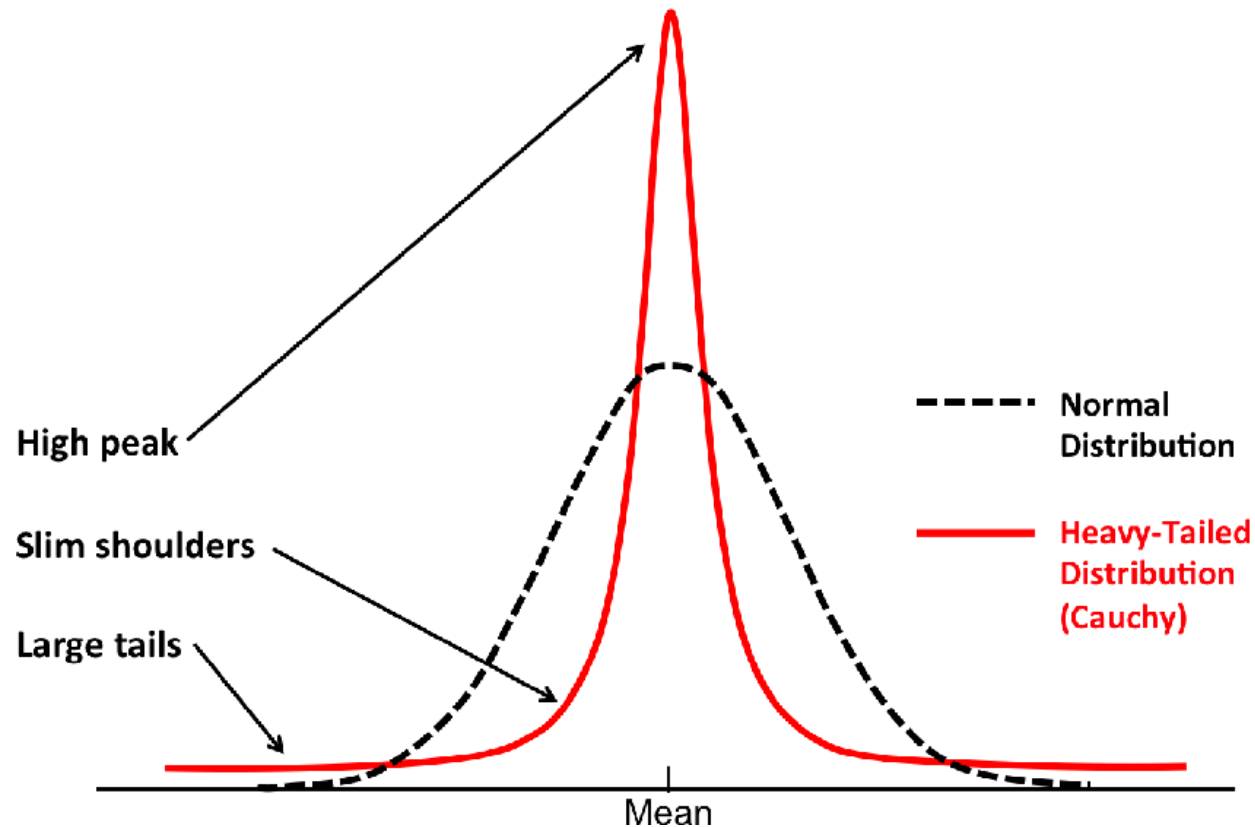
Specific form of linear regression

- So far we said linear regression is $\hat{y} = E_Y[p(y|x)] = E[p(y|x)] = \mathbf{w}^T \mathbf{x}$.
- But what do we use for $p(y|x)$?
- Standard linear regression uses a Gaussian $p(y|x) = N(y|\mathbf{w}^T \mathbf{x}, \sigma^2)$.
- Equivalent to $Y = \mathbf{w}^T \mathbf{x} + \epsilon$, with $\epsilon \sim N(0, \sigma^2)$.
- Which is equivalent to $Y - \mathbf{w}^T \mathbf{x} = \epsilon \sim N(0, \sigma^2)$.
- Alternate forms give "heavier tails" to the distribution of the "residual".



Aside: heavy-tailed distribution

- More of the mass lies further from the center of mass.
- Heavy-tailed noise models outliers better than a Gaussian.

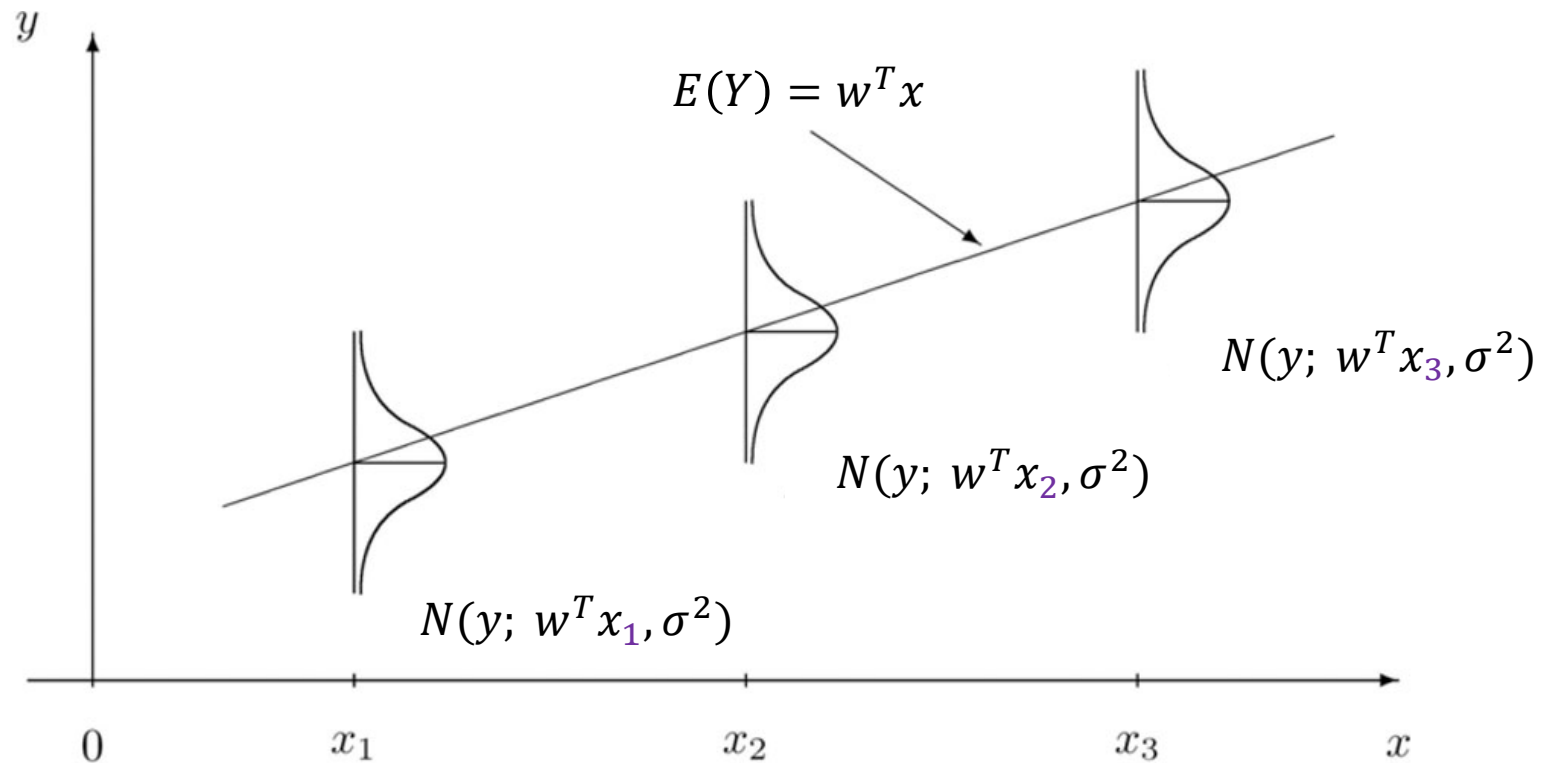


$$Y - w^T x = \epsilon \sim N(0, \sigma^2)$$

$$Y - w^T x = \epsilon \sim \text{Cauchy}(0, \sigma^2)$$

Gaussian linear regression, $p(y|x) = N(y|w^T x, \sigma^2)$

For every value $X = x$, the target variable, Y , takes on a Gaussian distribution with the same variance, σ^2 :



“Training” a Gaussian linear regression model

How will we fit the regression model, $p(y|x) = N(y|w^T x, \sigma^2)$?

$$\text{MLE: } \theta_{MLE} = (w_{MLE}, \sigma_{MLE}^2) = \arg \max_{(w, \sigma^2)} \log p(D = \{(x_i, y_i)_{i=1}^n\} | \theta)$$

$$= \arg \max_{(w, \sigma^2)} \sum_{i=1}^n \log p(y_i | x_i, \theta)$$

$$= \arg \max_{(w, \sigma^2)} \sum_{i=1}^n \log N(y_i | w^T x_i, \sigma^2)$$

$$= \arg \max_{(w, \sigma^2)} \sum_{i=1}^n \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - w^T x_i)^2$$

$$= \arg \max_{(w, \sigma^2)} n \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - w^T x_i)^2$$

If we assume σ^2 is known...

“Training” a Gaussian linear regression model

If we assume σ^2 is known...

$$(\mathbf{w}_{MLE}, \sigma_{MLE}^2) = \arg \max_{(\mathbf{w}, \sigma^2)} n \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

... then estimating \mathbf{w} above is the same as

$$\begin{aligned} &= \arg \min_{\mathbf{w}} \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2 \\ &= \arg \min_{\mathbf{w}} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \end{aligned}$$

“least squares” loss function!

$$p(y|x) = N(y|\mathbf{w}^T \mathbf{x}, \sigma^2)$$

"Training" a Gaussian linear regression model

$$w_{\text{MLE}} = \arg \min_w \sum_{i=1}^n (y_i - w^T x_i)^2$$

Lets re-write this loss in "vectorized" form:

First, define:

The image shows a handwritten derivation of the vectorized form of the loss function. The main equation is:

$$\begin{pmatrix} y_1 - w^T x_1 \\ \vdots \\ y_n - w^T x_n \end{pmatrix} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} - \begin{pmatrix} x_1^T w \\ \vdots \\ x_n^T w \end{pmatrix} = y - \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix} w = y - Aw$$

To the right of this equation, there is a smaller box containing the identity $w^T x_i = x_i^T w$. Below this box, the word "Scalar" is written with a checkmark above it. A blue arrow points from the "Scalar" label to the term $x_i^T w$ in the main equation.

"Training" a Gaussian linear regression model

$$w_{\text{MLE}} = \arg \min_w \sum_{i=1}^n (y_i - w^T x_i)^2$$

Lets re-write this loss in "vectorized" form:

First, define:

$$\begin{pmatrix} y_1 - w^T x_1 \\ \vdots \\ y_n - w^T x_n \end{pmatrix} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} - \begin{pmatrix} x_1^T w \\ \vdots \\ x_n^T w \end{pmatrix} = y - \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix} w = y - A w$$

$w^T x_i = x_i^T w$
Scalar

A is called
the "design
matrix"

"Training" a Gaussian linear regression model

$$\arg \min_w \sum_{i=1}^n (y_i - w^T x_i)^2$$

Lets re-write this loss in "vectorized" form:

First, define:

The image shows a handwritten derivation on a black background. The main equation is:
$$\begin{pmatrix} y_1 - w^T x_1 \\ \vdots \\ y_n - w^T x_n \end{pmatrix} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} - \begin{pmatrix} x_1^T w \\ \vdots \\ x_n^T w \end{pmatrix} = y - \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix} w = y - Aw$$
 To the right of this equation, there is a smaller box containing the equation $w^T x_i = x_i^T w$. Below this box, the word "Scalar" is written with a tilde symbol above it. A blue arrow points from the "Scalar" text to the $x_i^T w$ term in the main equation. A purple arrow points from the A matrix in the main equation to the text "A is called the 'design matrix'" located at the bottom right of the slide.

Then, we can re-write the loss as

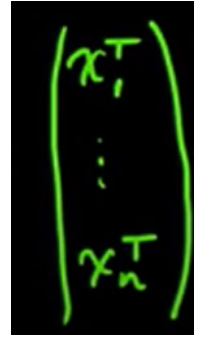
$$\begin{aligned} & \arg \min_w (y - Aw)^T (y - Aw) \\ &= \arg \min_w \|y - Aw\|_2^2 \end{aligned}$$

*A is called
the "design
matrix"*

"Training" a Gaussian linear regression model

So want to minimize

$$\begin{aligned}\mathcal{L} &= (\mathbf{y} - \mathbf{A}\mathbf{w})^T (\mathbf{y} - \mathbf{A}\mathbf{w}) && (\mathbf{y} \in \mathbb{R}^n, \mathbf{A} \in \mathbb{R}^{n \times d}, \mathbf{w} \in \mathbb{R}^{d \times 1}) \\ &= (\mathbf{y}^T - (\mathbf{A}\mathbf{w})^T)(\mathbf{y} - \mathbf{A}\mathbf{w}) \\ &= \mathbf{y}^T \mathbf{y} - \mathbf{w}^T \mathbf{A}^T \mathbf{y} - \mathbf{y}^T \mathbf{A}\mathbf{w} + \mathbf{w}^T \mathbf{A}^T \mathbf{A}\mathbf{w} \\ &= \mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T \mathbf{A}^T \mathbf{y} + \mathbf{w}^T \mathbf{A}^T \mathbf{A}\mathbf{w}\end{aligned}$$



\mathbf{A}

To minimize, we want to set $\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{0}$.

This is most easily achieved by using the rules of vector calculus, so let's do a quick refresher.

$$p(\mathbf{y}|\mathbf{x}) = N(\mathbf{y}|\mathbf{w}^T \mathbf{x}, \sigma^2)$$

Refresher on vector calculus

Some “rules” for taking gradients with respect to vectors.

- e.g., for vectors $\mathbf{a}, \mathbf{b} \in \mathbb{R}^{d \times 1}$, so that $\mathbf{a}^T \mathbf{b} \in \mathbb{R}$,

$$\frac{\partial(\mathbf{a}^T \mathbf{b})}{\partial a_j} = \frac{\partial(a_1 b_1 + a_2 b_2 + \cdots a_d b_d)}{\partial a_j} = b_j$$

Thus,

$$\frac{\partial(\mathbf{a}^T \mathbf{b})}{\partial \mathbf{a}} = \frac{\partial(a_1 b_1 + a_2 b_2 + \cdots a_d b_d)}{\partial \mathbf{a}} = \mathbf{b} \in \mathbb{R}^{d \times 1}$$

(not true for $\mathbf{a}\mathbf{b}^T$ which is a matrix, be careful!)

$$= \frac{\partial(\mathbf{b}^T \mathbf{a})}{\partial \mathbf{a}}$$

Refresher on vector calculus

For vector $\boldsymbol{x} \in \mathbb{R}^{d \times 1}$, and matrix $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$

$$\frac{\partial \boldsymbol{x}^T \boldsymbol{\Sigma} \boldsymbol{x}}{\partial \boldsymbol{x}} = (\boldsymbol{\Sigma} + \boldsymbol{\Sigma}^T) \boldsymbol{x}$$

Thus, if $\boldsymbol{\Sigma}$ is symmetric such that $\boldsymbol{\Sigma} = \boldsymbol{\Sigma}^T$ then

$$\frac{\partial \boldsymbol{x}^T \boldsymbol{\Sigma} \boldsymbol{x}}{\partial \boldsymbol{x}} = 2\boldsymbol{\Sigma} \boldsymbol{x}$$

(similar to the scalar version: $\frac{\partial (ax^2)}{\partial x} = 2ax$)

"Training" a Gaussian linear regression model

So want to minimize

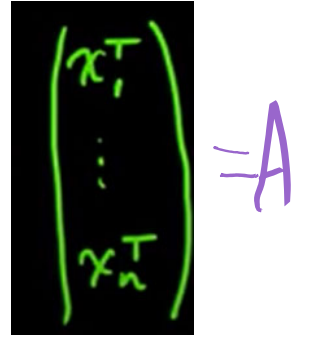
$$\mathcal{L} = (y - Aw)^T (y - Aw) \quad (y \in \mathbb{R}^n, A \in \mathbb{R}^{n \times d}, w \in \mathbb{R}^{d \times 1})$$

$$= (y^T - (Aw)^T)(y - Aw)$$

$$= y^T y - w^T A^T y - y^T Aw + w^T A^T Aw$$

$$= y^T y - 2w^T A^T y + w^T A^T Aw$$

To minimize, we want to set $\frac{\partial \mathcal{L}}{\partial w} = 0$.

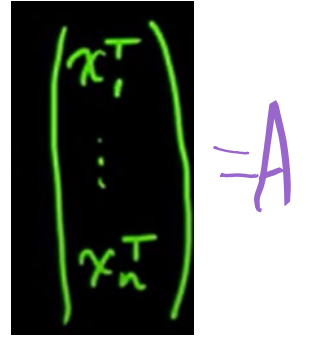


A handwritten diagram showing a vertical column of elements enclosed in large parentheses. The elements are labeled x_1^T at the top, a vertical ellipsis \vdots in the middle, and x_n^T at the bottom. To the right of the parentheses is a purple equals sign followed by the letter A .

"Training" a Gaussian linear regression model

So want to minimize

$$\mathcal{L} = (y - Aw)^T (y - Aw) \quad (y \in \mathbb{R}^n, A \in \mathbb{R}^{n \times d}, w \in \mathbb{R}^{d \times 1})$$



A handwritten diagram showing a vertical column of elements: x_1^T , a vertical ellipsis, and x_n^T . To the right of this column is a purple arrow pointing to the letter A .

$$= (y^T - (Aw)^T)(y - Aw)$$

$$= y^T y - w^T A^T y - y^T Aw + w^T A^T Aw$$

$$= y^T y - 2w^T A^T y + w^T A^T Aw$$

To minimize, we want to set $\frac{\partial \mathcal{L}}{\partial w} = 0$.

$$\nabla_w \mathcal{L} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial w_1} \\ \frac{\partial \mathcal{L}}{\partial w_2} \\ \vdots \\ \frac{\partial \mathcal{L}}{\partial w_d} \end{bmatrix} = \begin{bmatrix} -2A^T y + 2A^T A w \end{bmatrix}$$
$$= \underbrace{-A^T y + A^T A w}_{d \times 1}$$

"Training" a Gaussian linear regression model

So want to minimize

$$\begin{aligned}\mathcal{L} &= (y - Aw)^T (y - Aw) \\ &= (y^T - (Aw)^T)(y - Aw) \\ &= y^T y - w^T A^T y - y^T Aw + w^T A^T Aw \\ &= y^T y - 2w^T A^T y + w^T A^T Aw\end{aligned}$$

To minimize, we want to set $\frac{\partial \mathcal{L}}{\partial w} = 0$.

$$\begin{aligned}\nabla_w \mathcal{L} &= \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial w_1} \\ \frac{\partial \mathcal{L}}{\partial w_2} \\ \vdots \\ \frac{\partial \mathcal{L}}{\partial w_n} \end{bmatrix} = \begin{bmatrix} -2A^T y + 2A^T Aw \end{bmatrix} \\ &= \underbrace{-A^T y + A^T Aw}_{dx1}\end{aligned}$$

"left pseudo-inverse" of A
 $(A^T A)^{-1} A^T * A = I$

$$A^+ A = I$$

$$\begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix} = A$$

set to 0

rectangular
So no A^{-1} !

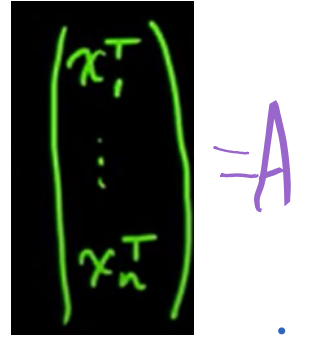
$$-A^T y + A^T Aw = 0$$

$$A^T Aw = A^T y$$

$$w = (A^T A)^{-1} A^T y$$

"Training" a Gaussian linear regression model

- We still need to check if the critical point, $w = (A^T A)^{-1} A^T y$ is minimum of the squared error loss.
- Recall $\nabla_w \mathcal{L} = -2A^T y + 2A^T A w$
- So Hessian matrix ($\nabla_w^2 \mathcal{L}$) is $2A^T A$. When is $A^T A$ PD?
- When the features are independent (when it has full rank).
- σ^2 from MLE as well is just the mean squared residual, $\sigma^2 = \frac{1}{N} \sum_i (y_i - w^T x)^2$.

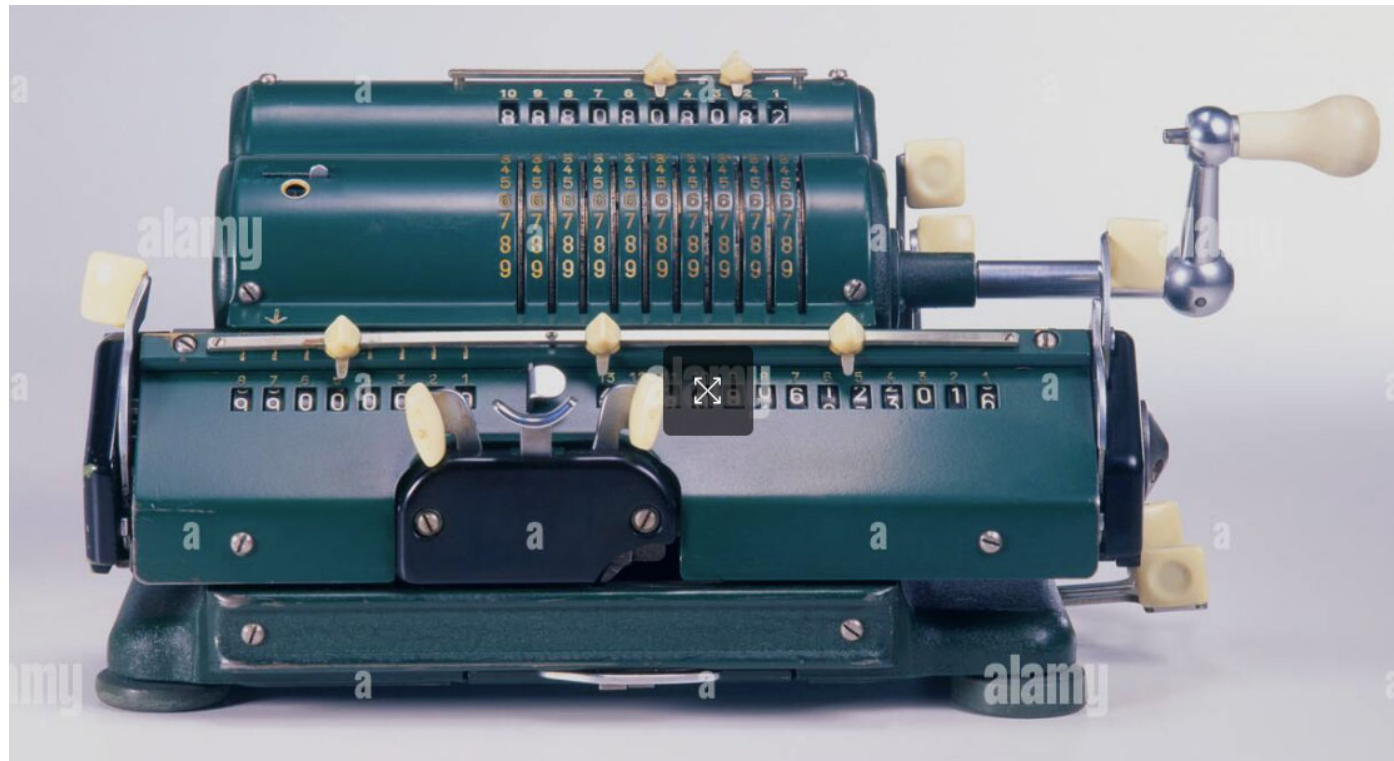
A handwritten diagram showing a vertical column of green parentheses containing the terms x_1^T , a vertical ellipsis, and x_n^T . To the right of the column is a purple equals sign followed by the letter A .

$$\nabla_w \mathcal{L} = \begin{bmatrix} \partial \mathcal{L} / \partial w_1 \\ \partial \mathcal{L} / \partial w_2 \\ \vdots \\ \partial \mathcal{L} / \partial w_n \end{bmatrix} = \begin{bmatrix} -2A^T y + 2A^T A w \end{bmatrix}$$
$$= \underbrace{-A^T y + A^T A w}_{dx_1}$$

$$\begin{aligned} -A^T y + A^T A w &= 0 \\ A^T A w &= A^T y \\ \boxed{w} &= \boxed{(A^T A)^{-1} A^T y} \end{aligned}$$

Regression in 1950s

Electromechanical desk "calculators" were used, and it could take up to 24 hours to receive the result from one regression.



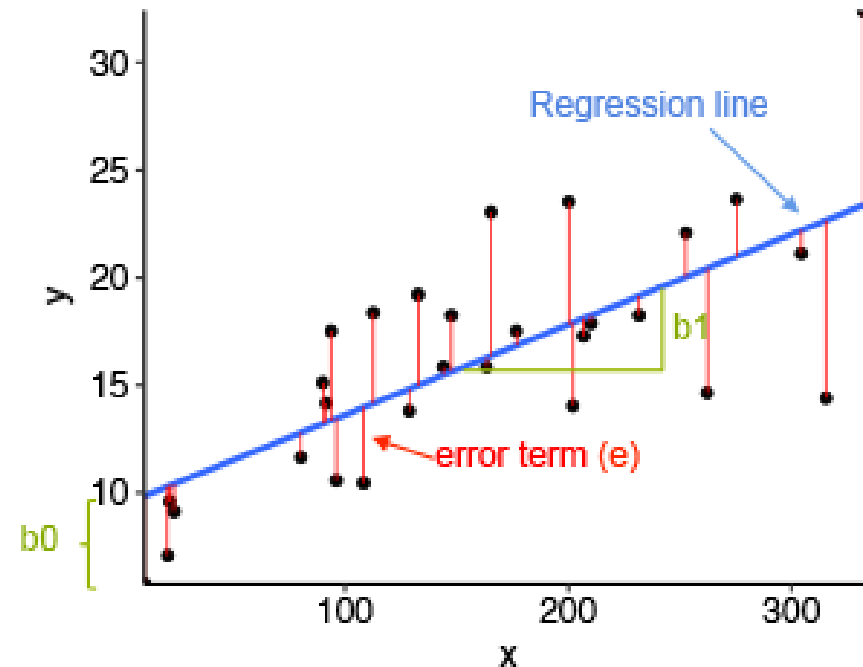
Geometric view of linear regression

- We're trying to predict all our training data labels correctly, such that $y_i = w^T x_i$ for all $i \in [1 \dots n]$.
- In vector form, this means we're looking for

$$\begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix} = A$$

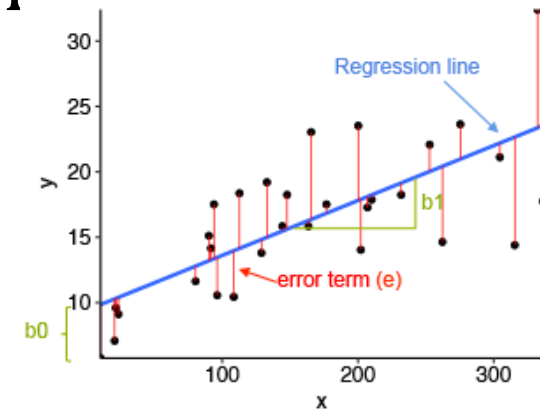
$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = A w = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix} w$$

Generally not possible because of noise; or incorrect model (e.g. all features).

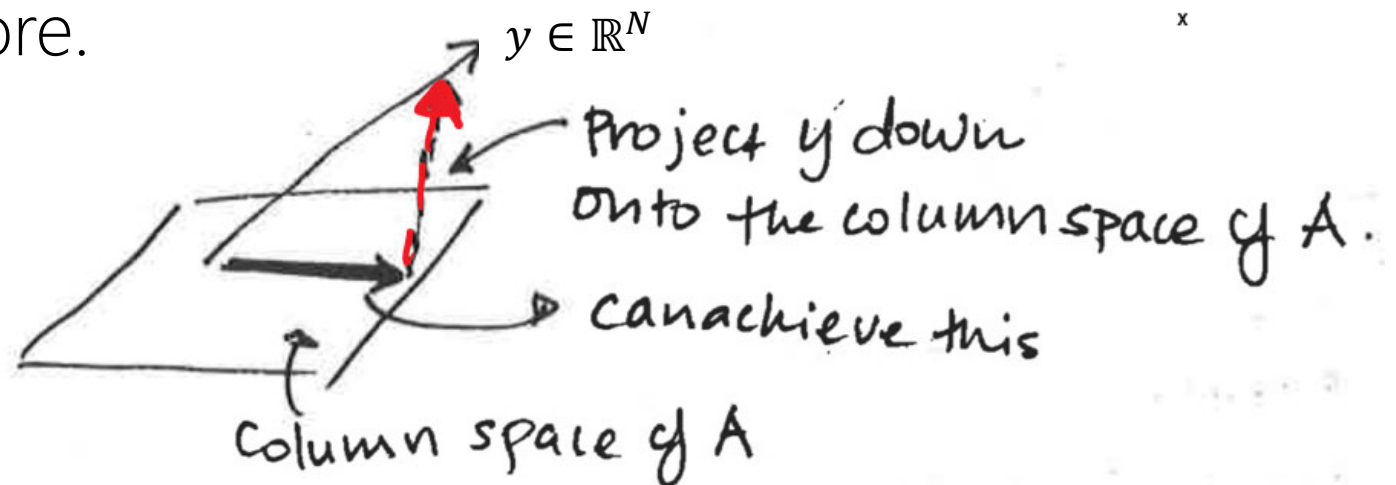


Geometric view of linear regression, $y = Aw$

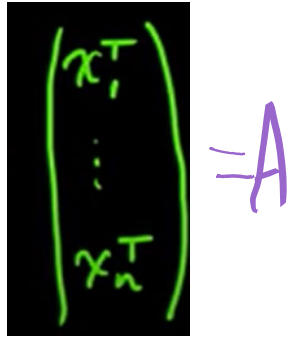
- So let's think about the **error vector** ($e \equiv y - \hat{y} = y - Aw, \in \mathbb{R}^{n \times 1}$).
- A "good" setting of w minimizes the length of e .
- The length is minimized when e lies \perp to column space of A
- Thus we seek w such that $A^T e = 0 = A^T (y - Aw)$.
- Thus $A^T y - A^T Aw = 0$.
- Thus $A^T y = A^T Aw$ (same as from MLE/least squares)!
- Thus $w = (A^T A)^{-1} A^T y$, as before.



$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} = Aw = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{bmatrix} w$$



Using basis expansions instead of x_i

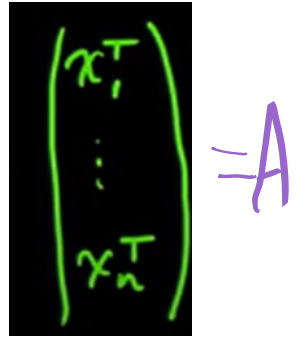
A handwritten matrix A is shown on a black background. The matrix is a column vector with two entries, x_1^T and x_2^T , separated by a vertical ellipsis. The matrix is enclosed in large parentheses. To the right of the matrix, there is a purple equals sign followed by the letter A .

- $\{x_j\} \rightarrow \{\Phi(x_j)\}$?
- Just define \mathbf{A} with $\Phi^T(x_j)$ instead of x_j because $\Phi(x)$ is fixed ahead of time, so its like someone just gave us different raw inputs x .

$$x = \Phi[x_1, x_2] = [1, x_1, x_2, x_1 x_2, x_1^2, x_2^2] \in \mathbb{R}$$

Polynomial expansion of order 2 (i.e. quadratic)

What can go wrong in linear regression

A hand-drawn diagram on a black background showing a vertical column of vectors. The top vector is labeled x_1^T , followed by a vertical ellipsis, and the bottom vector is labeled x_n^T . To the right of this column is a purple equals sign followed by the letter A .

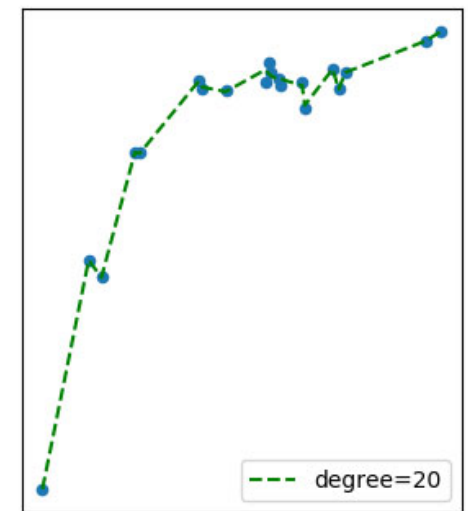
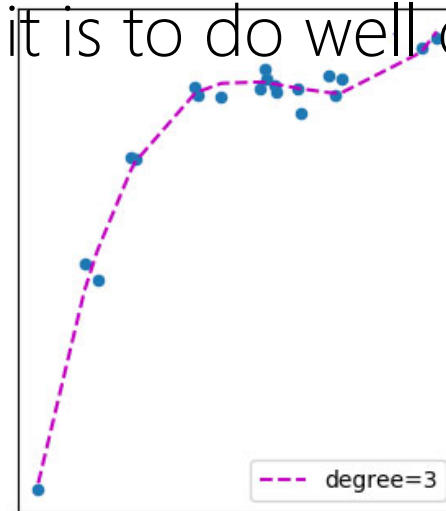
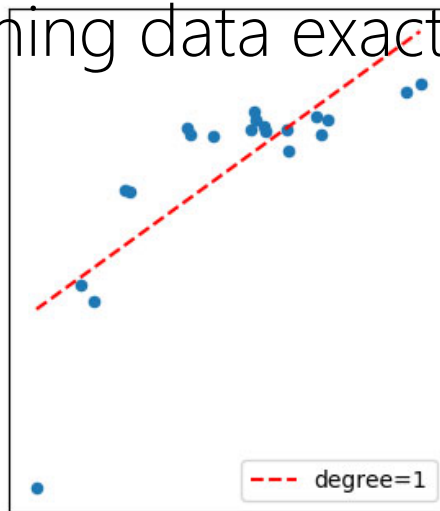
When can we invert $A^T A \in \mathbb{R}^{d \times d}$ (for d features)? ($A \in \mathbb{R}^{N \times d}$)

- Consider the decomposition we learned last class:
- $A^T A = \Sigma = Q D Q^T$ for diagonal D containing eigenvalues, and orthonormal Q .
- Then we had that $(A^T A)^{-1} = Q D^{-1} Q^T$. So if $A^T A$ has any zero eigenvalues we cannot invert it (there are ∞ equally good solutions w then).
- (Could use Moore-Penrose pseudo inverse.)
- This degeneracy occurs when the features in A , for the given data $\{x_j\}$ are linearly dependent.
- e.g., when $d > n$ (then $\text{rank}(A^T A) < n$, but needs to be D , i.e. "full rank" to avoid zero eigenvalues).
- What about when $d = n$?

What can go wrong in linear regression

As we add higher and higher order polynomials, a few things happen:

1. # features, d gets bigger and bigger.
2. For $d \geq n$ can perfectly fit any data (i.e., polynomials are a complete basis).
3. Even when we don't perfectly fit the training data, we are still in danger of overfitting (worse prediction on test set). Our goal is not to fit a line through the training data exactly, it is to do well on unseen test cases!



What can go wrong in linear regression

Two main categories of fixes:

1. Remove features until the problem is well behaved.
2. Leave the features as they are, but add constraints to the system to “tighten it up” (aka “regularization”).
 - 1) Is called “feature selection”, e.g. “forward selection”, “backward selection”, etc.

(Moore-Penrose inverse is not a general fix for ML models, only works for linear regression.)

A thought experiment

Consider:

- Use MLE on data, $D = \{x_i, y_i\}$ to get $\hat{y}_i = p_{\theta}(y|\Phi(x))$ using linear regression.
- Assume an abundance of data (millions of data points), and only 100 parameters.
- Suppose get accuracy $\mp \$1000$ of sale price when applying to held out part of our data.
- Can we assume this model will get $\mp \$1000$ on any test set that may come in the future?

Actual vs. predicted sale price of house

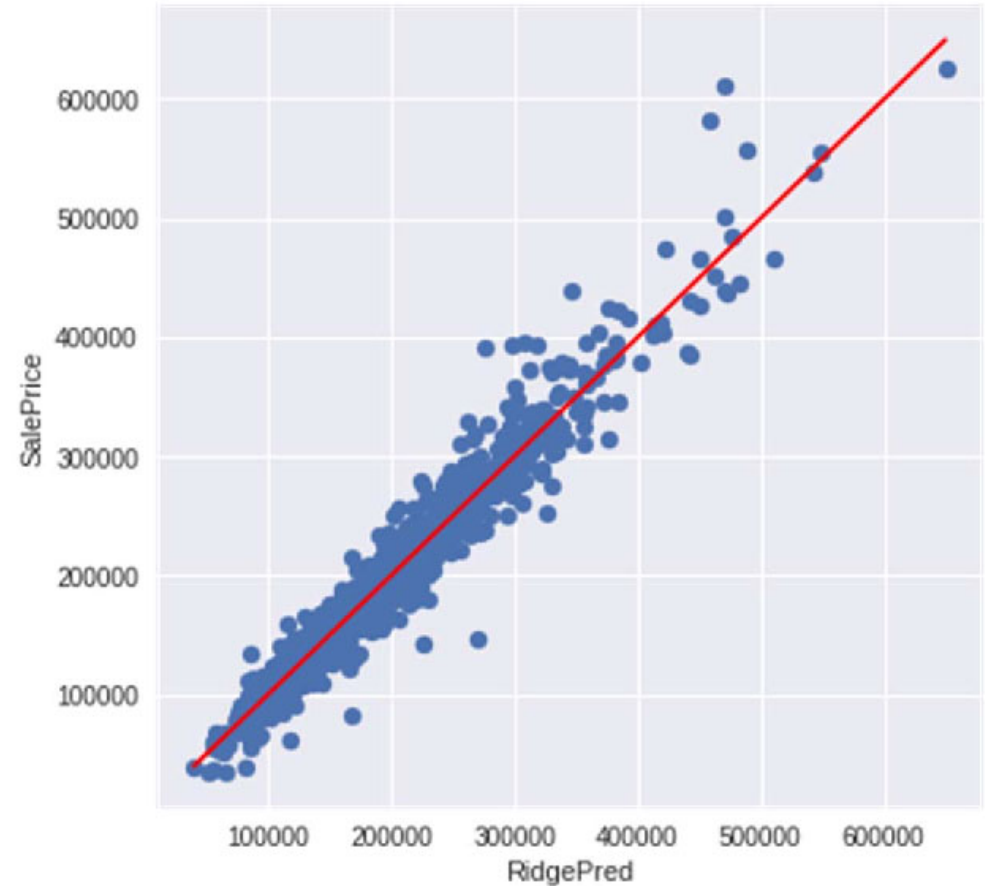


Fig. 4 Ridge Prediction for Training Data.

Causation vs correlation

Breakingviews

Zillow's failed house flipping

Reuters

WSJ NOV. 2021 : *“The company expects to record losses of more than \$500 million from home-flipping by the end of this year and is laying off a quarter of its staff.”*

Actual vs. predicted sale price of house

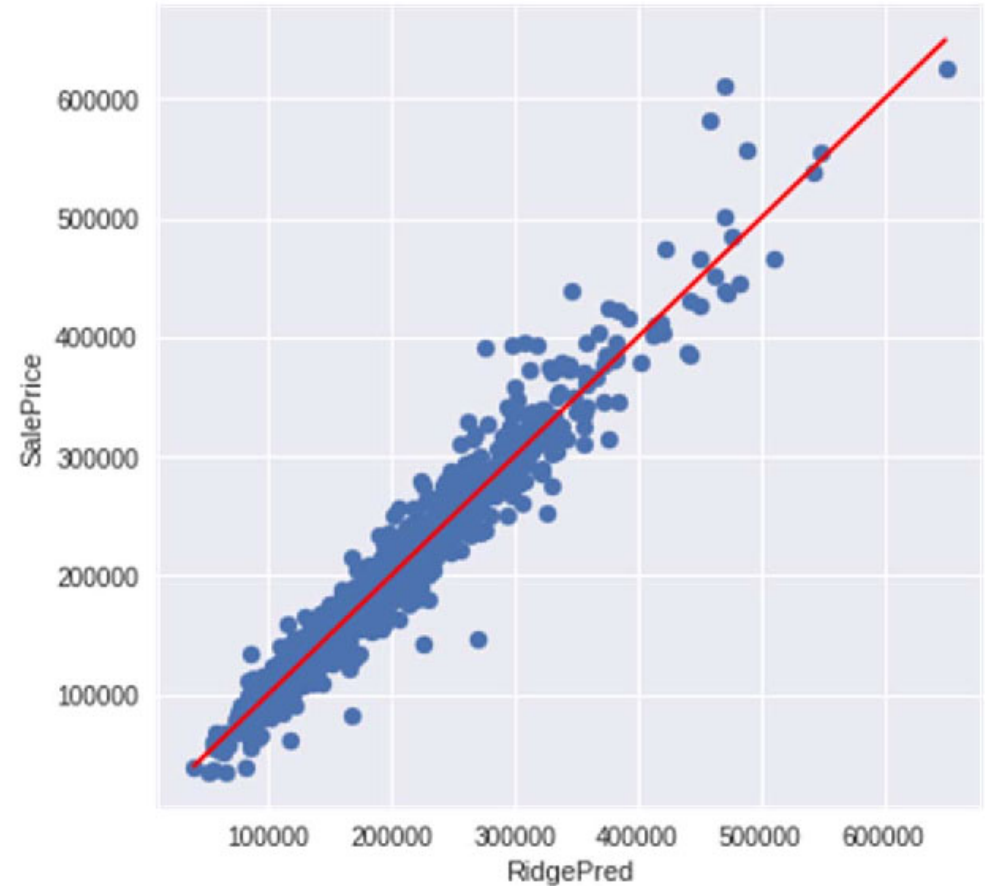


Fig. 4 Ridge Prediction for Training Data.