



Linear Regression

Machine Learning – CSE546

Kevin Jamieson

University of Washington

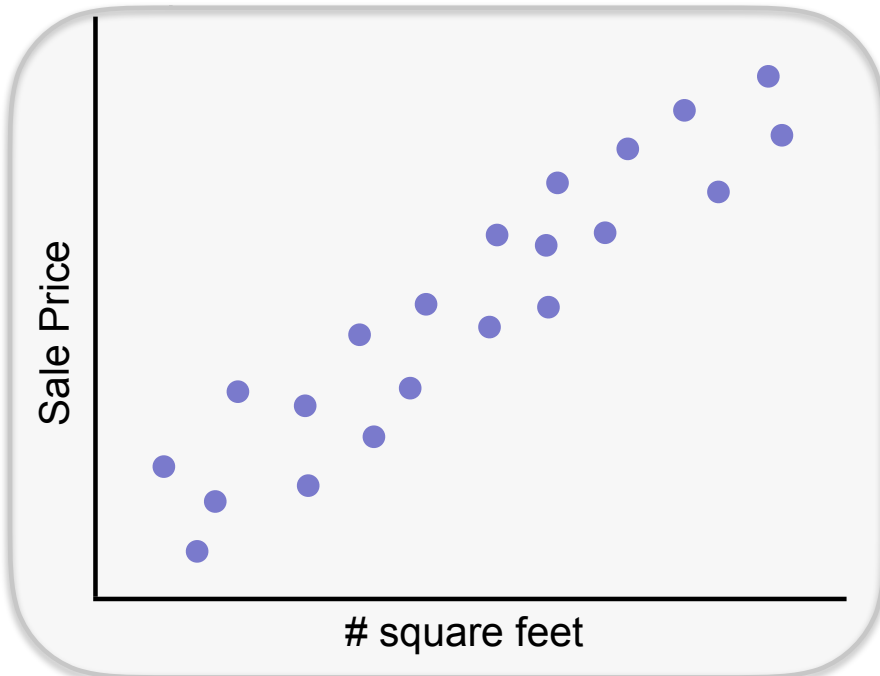
Oct 2, 2018

The regression problem

Given past sales data on [zillow.com](https://www.zillow.com), predict:

$y =$ **House sale price** *from*

$x =$ **{# sq. ft., zip code, date of sale, etc.}**



Training Data:

$$\{(x_i, y_i)\}_{i=1}^n$$

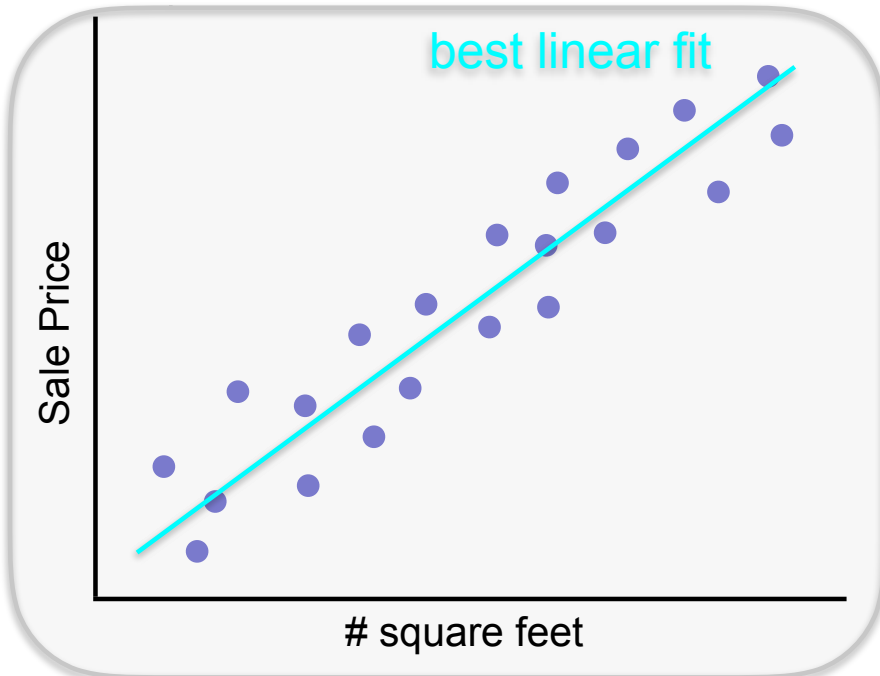
$$x_i \in \mathbb{R}^d$$
$$y_i \in \mathbb{R}$$

The regression problem

Given past sales data on [zillow.com](https://www.zillow.com), predict:

$y =$ **House sale price** *from*

$x =$ **{# sq. ft., zip code, date of sale, etc.}**



Training Data:

$$\{(x_i, y_i)\}_{i=1}^n$$

$$x_i \in \mathbb{R}^d \\ y_i \in \mathbb{R}$$

Hypothesis: linear

$$y_i \approx x_i^T w$$

Loss: least squares

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

The regression problem in matrix notation

$$\begin{aligned}\hat{w}_{LS} &= \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2 \\ &= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)\end{aligned}$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix}$$

The regression problem in matrix notation

$$\begin{aligned}\hat{w}_{LS} &= \arg \min_w \|\mathbf{y} - \mathbf{X}w\|_2^2 \\ &= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)\end{aligned}$$

$$\nabla_w(\cdot) = -2X^T(y - Xw) = 0$$
$$\boxed{X^T y = X^T X w}$$

If $(X^T X)^{-1}$ exists then $\hat{w} = (X^T X)^{-1} X^T y$

The regression problem in matrix notation

$$\begin{aligned}\hat{w}_{LS} &= \arg \min_w \|\mathbf{y} - \mathbf{X}w\|_2^2 \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$

What about an offset?

$$\begin{aligned}\hat{w}_{LS}, \hat{b}_{LS} &= \arg \min_{w,b} \sum_{i=1}^n (y_i - (x_i^T w + b))^2 \\ &= \arg \min_{w,b} \|\mathbf{y} - (\mathbf{X}w + \mathbf{1}b)\|_2^2\end{aligned}$$

Dealing with an offset

$$\mathbf{1} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

$$\hat{w}_{LS}, \hat{b}_{LS} = \arg \min_{w, b} \|\mathbf{y} - (\mathbf{X}w + \mathbf{1}b)\|_2^2$$

$$\nabla_w = 2(\mathbf{X}^T)(\mathbf{y} - (\mathbf{X}w + \mathbf{1}b)) = 0$$

$$\mathbf{X}^T \mathbf{y} = \mathbf{X}^T \mathbf{X} w + \mathbf{X}^T \mathbf{1} b$$

$$\nabla_b = 2 \mathbf{1}^T (\mathbf{y} - (\mathbf{X}w + \mathbf{1}b)) = 0$$

$$\mathbf{1}^T \mathbf{1} = n$$

$$\mathbf{1}^T \mathbf{y} = \mathbf{1}^T \mathbf{X} w + n b$$

$$\mathbf{1}^T \mathbf{y} = \sum_{i=1}^n y_i$$

$$b = \frac{1}{n} \sum y_i - \frac{1}{n} \mathbf{1}^T \mathbf{X} w$$

Dealing with an offset

$$\hat{w}_{LS}, \hat{b}_{LS} = \arg \min_{w, b} \|\mathbf{y} - (\mathbf{X}w + \mathbf{1}b)\|_2^2$$

$$\mathbf{X}^T \mathbf{X} \hat{w}_{LS} + \hat{b}_{LS} \mathbf{X}^T \mathbf{1} = \mathbf{X}^T \mathbf{y}$$

$$\mathbf{1}^T \mathbf{X} \hat{w}_{LS} + \hat{b}_{LS} \mathbf{1}^T \mathbf{1} = \mathbf{1}^T \mathbf{y} \quad \left(\mathbf{1}^T \mathbf{X} w \right)^T = w^T \underbrace{\mathbf{X}^T \mathbf{1}}_{=0} = 0$$

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i$$

$$0 = \sum_{i=1}^n (x_i - \mu) = \left(\sum x_i \right) - n\mu = 0$$

If $\mathbf{X}^T \mathbf{1} = 0$ (i.e., if each feature is mean-zero) then

$$\hat{w}_{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

$$\begin{aligned} \tilde{x}_i &= x_i - \mu \\ \tilde{\mathbf{X}} &= \begin{bmatrix} \tilde{x}_1^T \\ \vdots \\ \tilde{x}_n^T \end{bmatrix} \end{aligned} \rightarrow \hat{w}$$

$$\hat{b}_{LS} = \frac{1}{n} \sum_{i=1}^n y_i \quad \text{Given a new } x, \text{ predict } \hat{w}^T (x - \mu) + \hat{b}$$

The regression problem in matrix notation

$$\begin{aligned}\hat{w}_{LS} &= \arg \min_w \|\mathbf{y} - \mathbf{X}w\|_2^2 \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$

But why least squares?

Consider $y_i = \underbrace{x_i^T w + \epsilon_i}$ where $\epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$

$$P(y|x, w, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - x^T w)^2}{2\sigma^2}\right)$$

Maximizing log-likelihood

Maximize:

$$\log P(\mathcal{D}|w, \sigma) = \log\left(\frac{1}{\sqrt{2\pi\sigma}}\right)^n \prod_{i=1}^n e^{-\frac{(y_i - x_i^T w)^2}{2\sigma^2}}$$

MLE is LS under linear model

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

$$\hat{w}_{MLE} = \arg \max_w P(\mathcal{D}|w, \sigma)$$

$$\text{if } y_i = x_i^T w + \epsilon_i \quad \text{and} \quad \epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$$

$$\hat{w}_{LS} = \hat{w}_{MLE} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

Analysis of error

$$\underline{\mathbf{Y}} = \mathbf{X}\mathbf{w} + \epsilon$$

if $y_i = x_i^T \mathbf{w} + \epsilon_i$ and $\epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$

$$\hat{\mathbf{w}}_{MLE} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

$$= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}\mathbf{w} + \boldsymbol{\epsilon})$$

$$= \mathbf{w} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon}$$

$$\mathbb{E}[\hat{\mathbf{w}}] = \mathbf{w}$$

$$\mathbb{E}[(\hat{\mathbf{w}} - \mathbf{w})(\hat{\mathbf{w}} - \mathbf{w})^T] = \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon} \boldsymbol{\epsilon}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}]$$

$$= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}[\boldsymbol{\epsilon} \boldsymbol{\epsilon}^T] \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}$$

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

$$\hat{\mathbf{w}} \sim \mathcal{N}(\mathbf{w}, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1})$$

$$\mathbb{E}[\boldsymbol{\epsilon} \boldsymbol{\epsilon}^T] = \sigma^2 \mathbf{I}$$

Analysis of error

$$\mathbf{Y} = \mathbf{X}w + \epsilon$$

if $y_i = x_i^T w + \epsilon_i$ and $\epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$

$$\begin{aligned}\hat{w}_{MLE} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}w + \epsilon) \\ &= w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon\end{aligned}$$

$$\text{Cov}(\hat{w}_{MLE}) = \mathbb{E}[(\hat{w} - \mathbb{E}[\hat{w}])(\hat{w} - \mathbb{E}[\hat{w}])^T] = (\mathbf{X}^T \mathbf{X})^{-1}$$

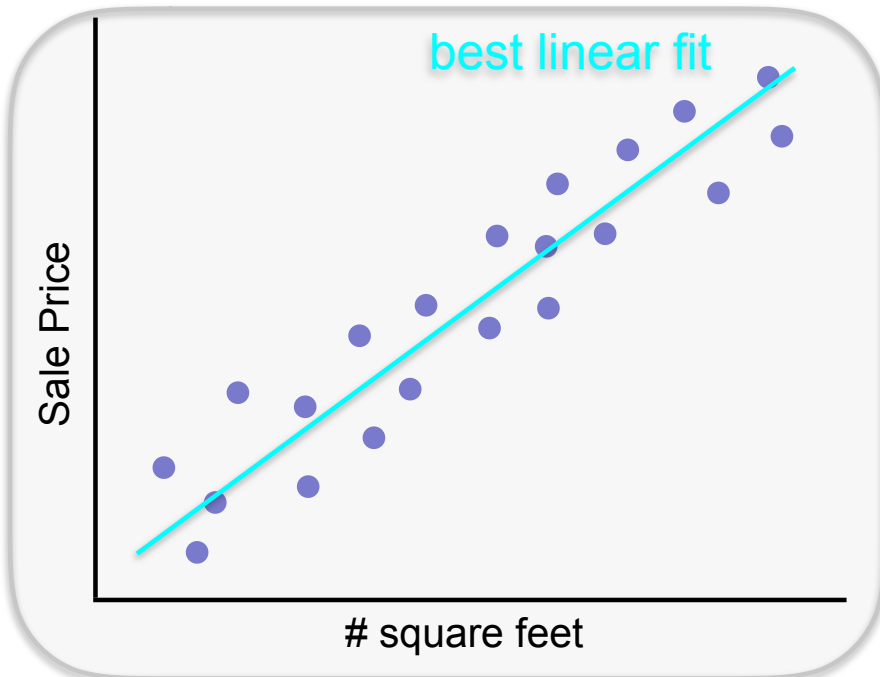
$$\hat{w}_{MLE} \sim \mathcal{N}(w, (\mathbf{X}^T \mathbf{X})^{-1})$$

The regression problem

Given past sales data on [zillow.com](https://www.zillow.com), predict:

$y =$ **House sale price** *from*

$x =$ **{# sq. ft., zip code, date of sale, etc.}**



Training Data:

$$\{(x_i, y_i)\}_{i=1}^n$$

$$x_i \in \mathbb{R}^d$$
$$y_i \in \mathbb{R}$$

Hypothesis: linear

$$y_i \approx x_i^T w$$

Loss: least squares

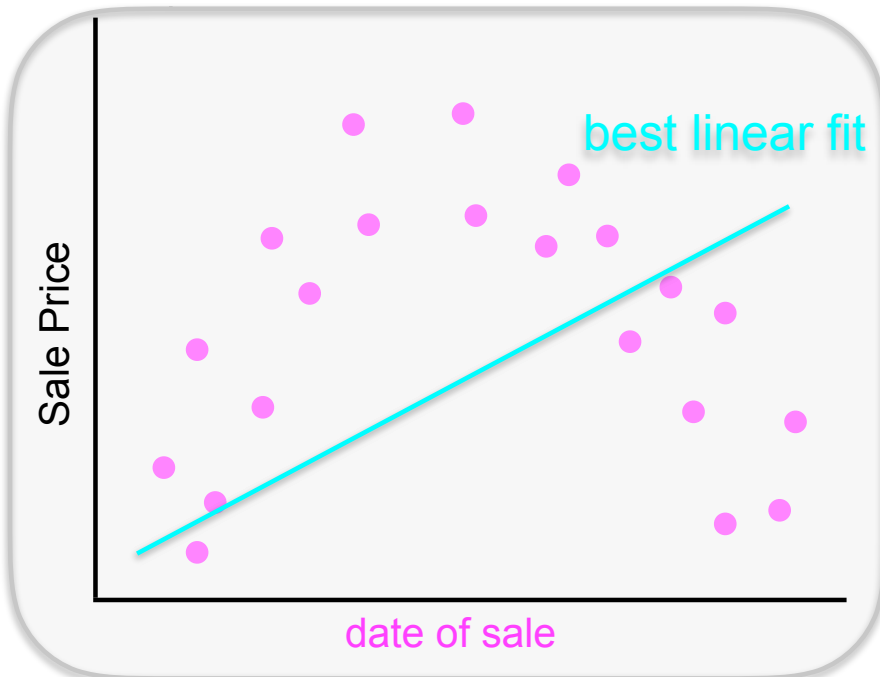
$$\min_w \sum_{i=1}^n (y_i - x_i^T w)^2$$

The regression problem

Given past sales data on [zillow.com](https://www.zillow.com), predict:

$y =$ **House sale price** *from*

$x =$ {# sq. ft., zip code, date of sale, etc.}



Training Data:

$$\{(x_i, y_i)\}_{i=1}^n$$

$$x_i \in \mathbb{R}^d$$
$$y_i \in \mathbb{R}$$

Hypothesis: linear

$$y_i \approx x_i^T w$$

Loss: least squares

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

The regression problem

Training Data: $x_i \in \mathbb{R}^d$
 $y_i \in \mathbb{R}$
 $\{(x_i, y_i)\}_{i=1}^n$

Transformed data:

Hypothesis: linear

$$y_i \approx x_i^T w$$

Loss: least squares

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

The regression problem

Training Data: $x_i \in \mathbb{R}^d$
 $y_i \in \mathbb{R}$
 $\{(x_i, y_i)\}_{i=1}^n$

Hypothesis: linear

$$y_i \approx x_i^T w$$

Loss: least squares

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

Transformed data:

$h : \mathbb{R}^d \rightarrow \mathbb{R}^p$ maps original features to a rich, possibly high-dimensional space

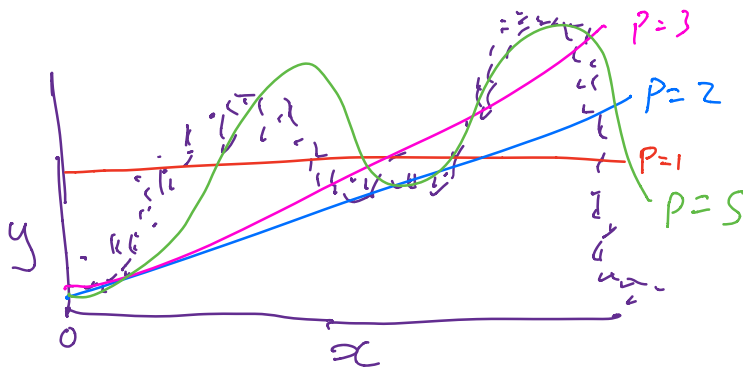
$$\text{in } d=1: h(x) = \begin{bmatrix} h_1(x) \\ h_2(x) \\ \vdots \\ h_p(x) \end{bmatrix} = \begin{bmatrix} x \\ x^2 \\ \vdots \\ x^p \end{bmatrix}$$

for $d>1$, generate $\{u_j\}_{j=1}^p \subset \mathbb{R}^d$

$$h_j(x) = \frac{1}{1 + \exp(u_j^T x)}$$

$$h_j(x) = (u_j^T x)^2$$

$$h_j(x) = \cos(u_j^T x)$$



$$\hat{f}(x) = \sum_{k=0}^{p-1} \underbrace{f^{(k)}(0)}_{\text{green box}} \frac{1}{k!} x^k$$

$$h(x) = \begin{bmatrix} h_0(x) \\ \vdots \\ h_{p-1}(x) \end{bmatrix}$$

Given x_i , predict $\tilde{w}^T h(x)$

$$h_k(x) = x^k \frac{1}{k!}$$

$$x_i \mapsto h(x_i) \quad \tilde{X} = \begin{bmatrix} h(x_1)^T \\ \vdots \\ h(x_n)^T \end{bmatrix}$$

$$\tilde{w} = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T y$$

The regression problem

Training Data: $x_i \in \mathbb{R}^d$
 $y_i \in \mathbb{R}$
 $\{(x_i, y_i)\}_{i=1}^n$

Hypothesis: linear

$$y_i \approx x_i^T w$$

Loss: least squares

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

Transformed data:

$$h(x) = \begin{bmatrix} h_1(x) \\ h_2(x) \\ \vdots \\ h_p(x) \end{bmatrix}$$

Hypothesis: linear

$$y_i \approx h(x_i)^T w \quad w \in \mathbb{R}^p$$

Loss: least squares

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

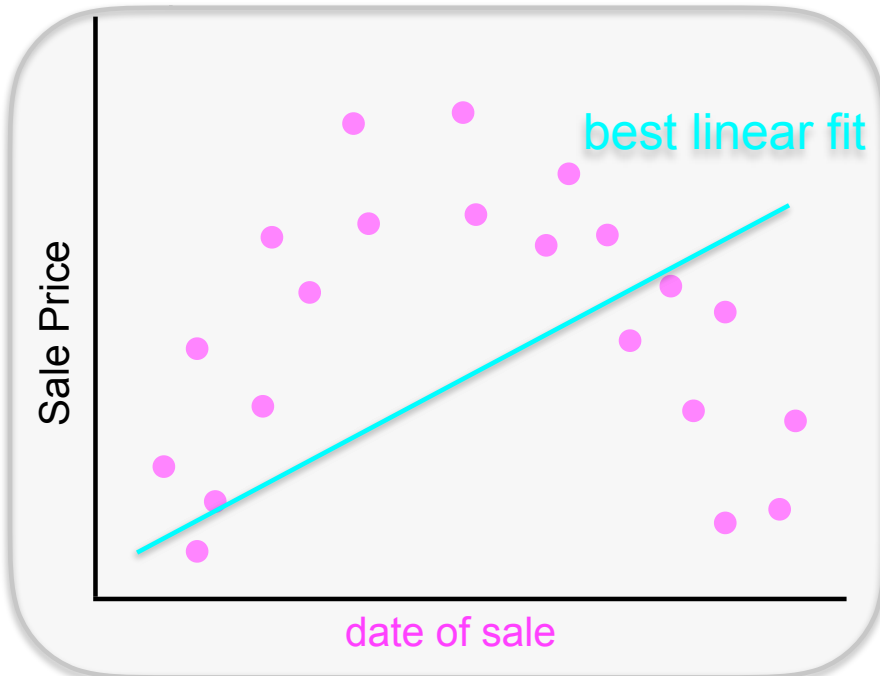
The regression problem

Training Data:

$$\{(x_i, y_i)\}_{i=1}^n \quad \begin{array}{l} x_i \in \mathbb{R}^d \\ y_i \in \mathbb{R} \end{array}$$

Transformed data:

$$h(x) = \begin{bmatrix} h_1(x) \\ h_2(x) \\ \vdots \\ h_p(x) \end{bmatrix}$$



Hypothesis: linear

$$y_i \approx h(x_i)^T w \quad w \in \mathbb{R}^p$$

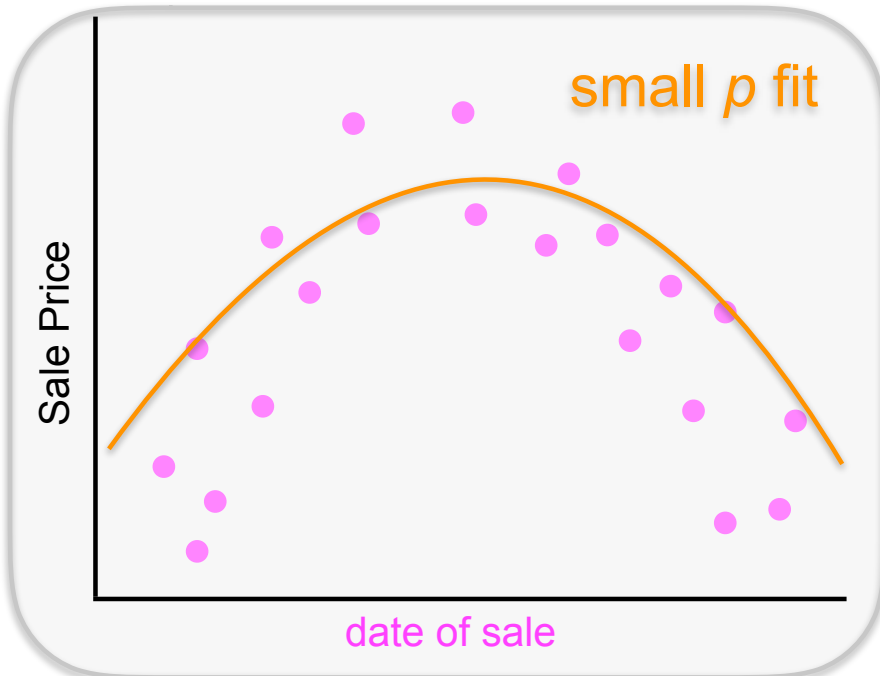
Loss: least squares

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

The regression problem

Training Data: $x_i \in \mathbb{R}^d$
 $y_i \in \mathbb{R}$
 $\{(x_i, y_i)\}_{i=1}^n$

Transformed data:

$$h(x) = \begin{bmatrix} h_1(x) \\ h_2(x) \\ \vdots \\ h_p(x) \end{bmatrix}$$


Hypothesis: linear

$$y_i \approx h(x_i)^T w \quad w \in \mathbb{R}^p$$

Loss: least squares

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

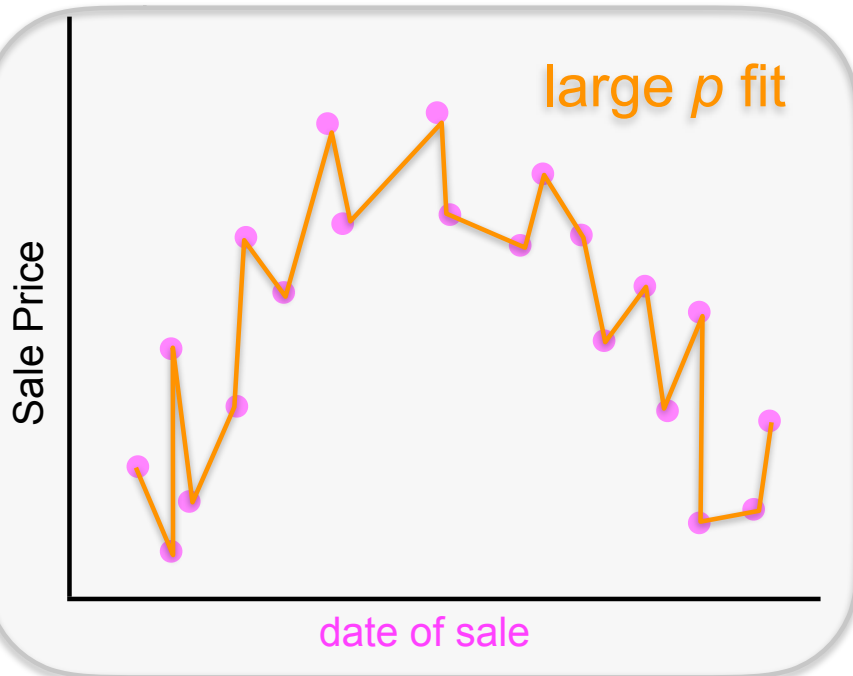
The regression problem

Training Data:

$$\{(x_i, y_i)\}_{i=1}^n \quad \begin{array}{l} x_i \in \mathbb{R}^d \\ y_i \in \mathbb{R} \end{array}$$

Transformed data:

$$h(x) = \begin{bmatrix} h_1(x) \\ h_2(x) \\ \vdots \\ h_p(x) \end{bmatrix}$$



Hypothesis: linear

$$y_i \approx h(x_i)^T w \quad w \in \mathbb{R}^p$$

Loss: least squares

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

What's going on here?



Bias-Variance Tradeoff

Machine Learning – CSE546

Kevin Jamieson

University of Washington

Oct 5, 2018

Statistical Learning

$$P_{XY}(X = x, Y = y)$$

Goal: Predict Y given X

Find function η that minimizes

$$\mathbb{E}_{XY}[(Y - \eta(X))^2] = \mathbb{E}_x \left[\mathbb{E}_{Y|X} [(Y - \eta(x))^2 | X=x] \right]$$

$$\eta(x) = \underset{c}{\operatorname{argmin}} \mathbb{E}_{Y|X} [(Y - c)^2 | X=x]$$

$$\begin{aligned} \frac{d}{dc} \mathbb{E}_{Y|X} [(Y - c)^2 | X=x] &= \mathbb{E} [2(Y - c) | X=x] = 0 \\ &= 2 \mathbb{E}[Y | X=x] - 2c = 0 \end{aligned}$$

$$c = \eta(x) = \mathbb{E}[Y | X=x]$$

Statistical Learning

$$P_{XY}(X = x, Y = y)$$

Goal: Predict Y given X

Find function η that minimizes

$$\mathbb{E}_{XY}[(Y - \eta(X))^2] = \mathbb{E}_X \left[\mathbb{E}_{Y|X}[(Y - \eta(x))^2 | X = x] \right]$$

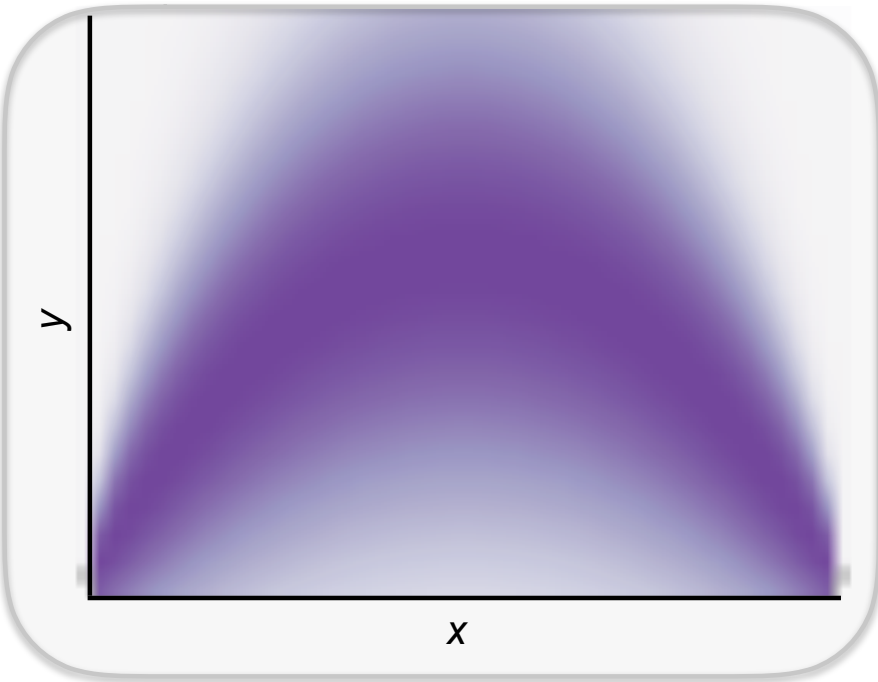
$$\eta(x) = \arg \min_c \mathbb{E}_{Y|X}[(Y - c)^2 | X = x] = \mathbb{E}_{Y|X}[Y | X = x]$$

Under LS loss, optimal predictor: $\eta(x) = \mathbb{E}_{Y|X}[Y | X = x]$

Statistical Learning

$$\mathbb{E}_{XY}[(Y - \eta(X))^2]$$

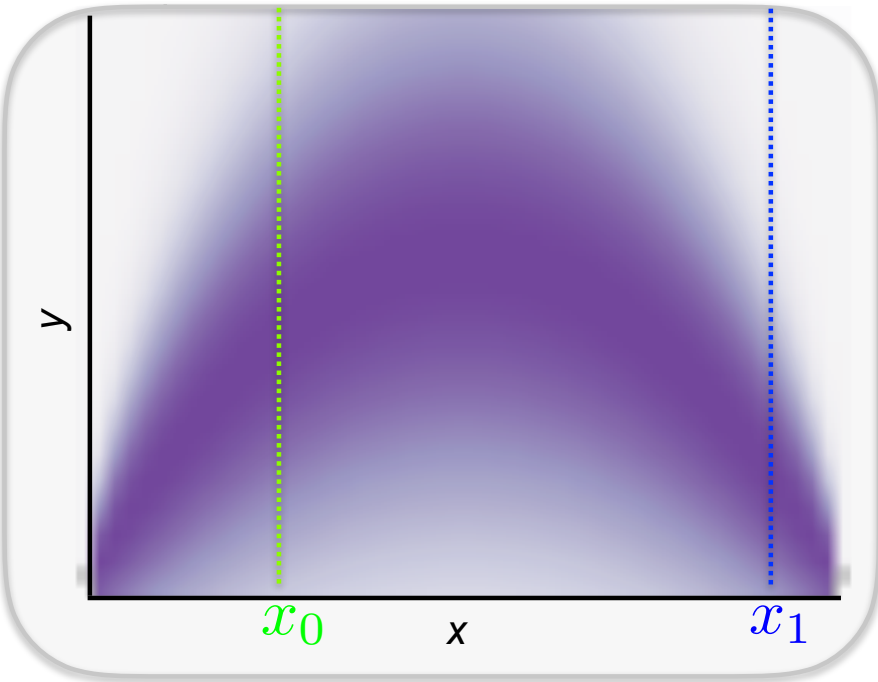
$$P_{XY}(X = x, Y = y)$$



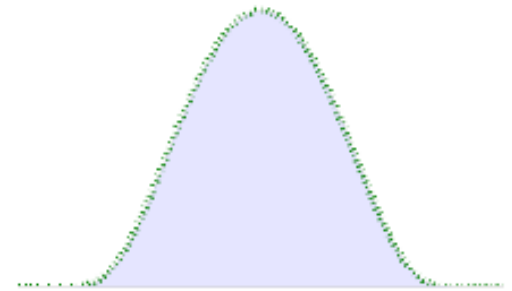
Statistical Learning

$$\mathbb{E}_{XY}[(Y - \eta(X))^2]$$

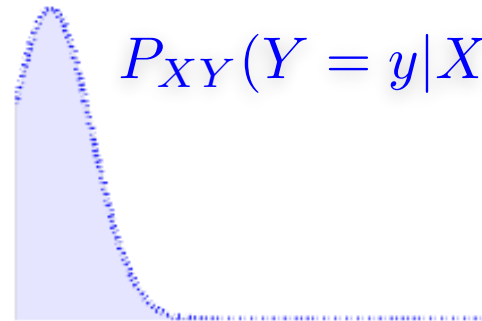
$$P_{XY}(X = x, Y = y)$$



$$P_{XY}(Y = y|X = x_0)$$



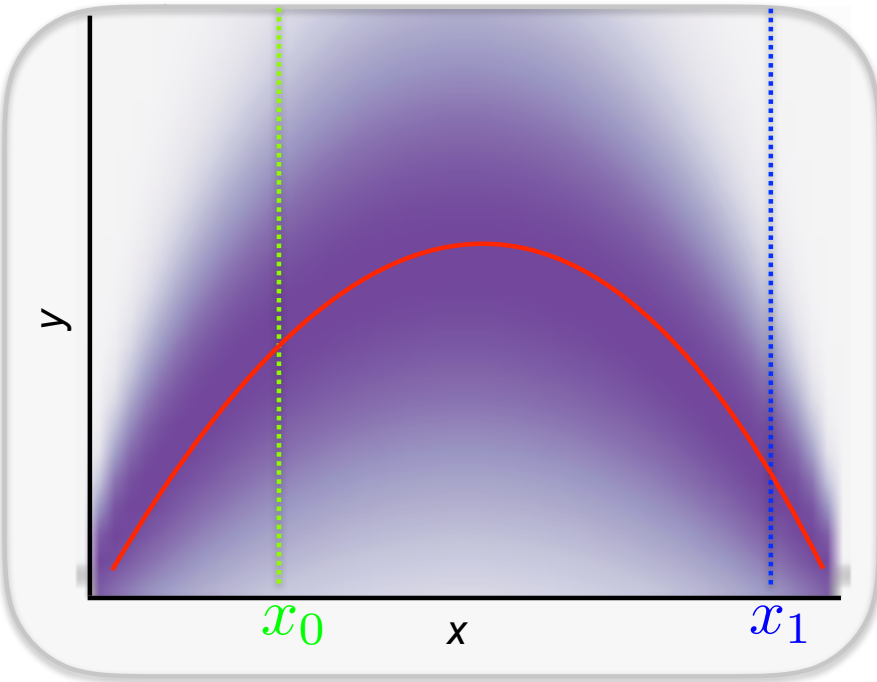
$$P_{XY}(Y = y|X = x_1)$$



Statistical Learning

$$\mathbb{E}_{XY}[(Y - \eta(X))^2]$$

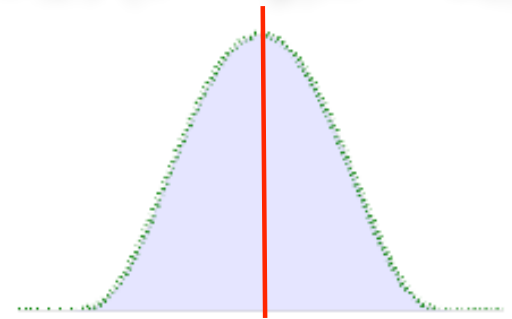
$$P_{XY}(X = x, Y = y)$$



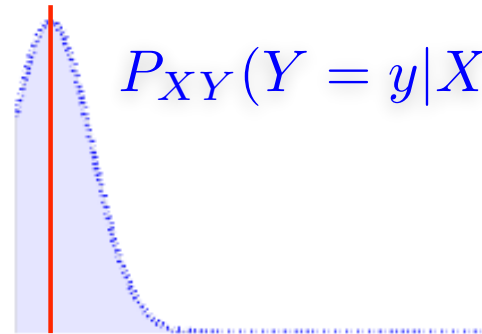
Ideally, we want to find:

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

$$P_{XY}(Y = y|X = x_0)$$



$$P_{XY}(Y = y|X = x_1)$$

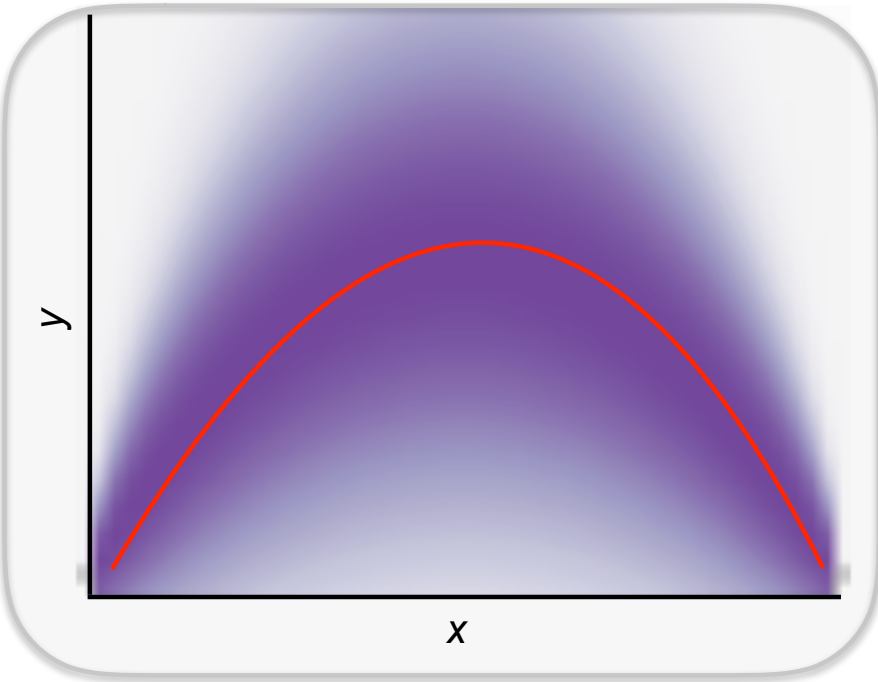


Statistical Learning

$$P_{XY}(X = x, Y = y)$$

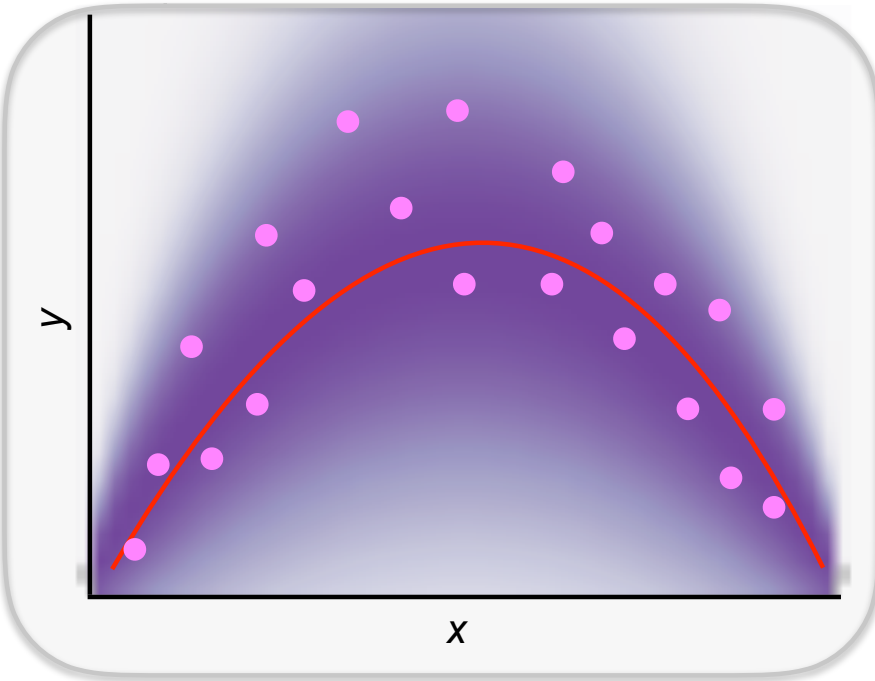
Ideally, we want to find:

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$



Statistical Learning

$$P_{XY}(X = x, Y = y)$$



Ideally, we want to find:

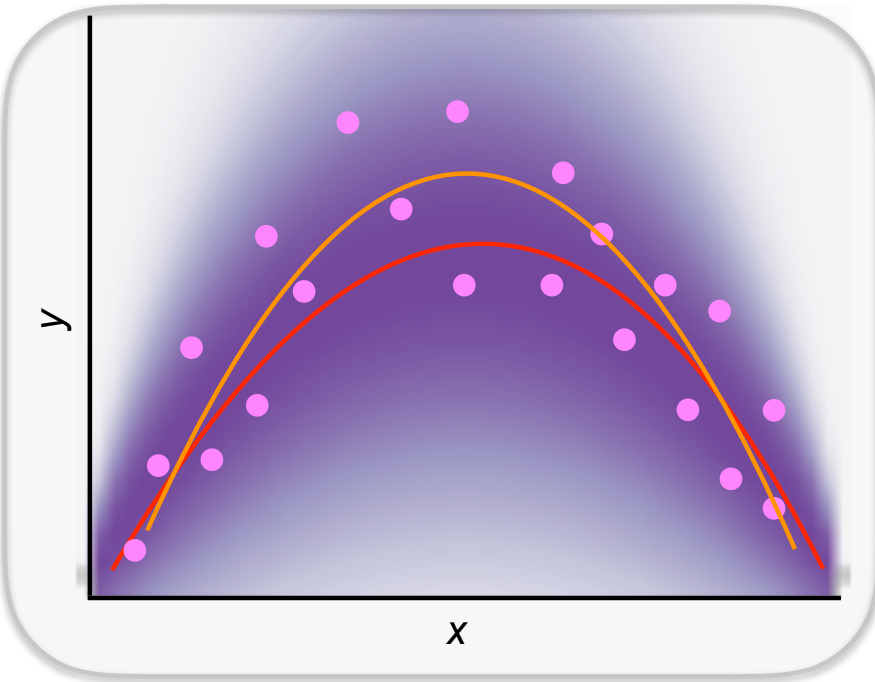
$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

But we only have samples:

$$(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY} \quad \text{for } i = 1, \dots, n$$

Statistical Learning

$$P_{XY}(X = x, Y = y)$$



Ideally, we want to find:

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

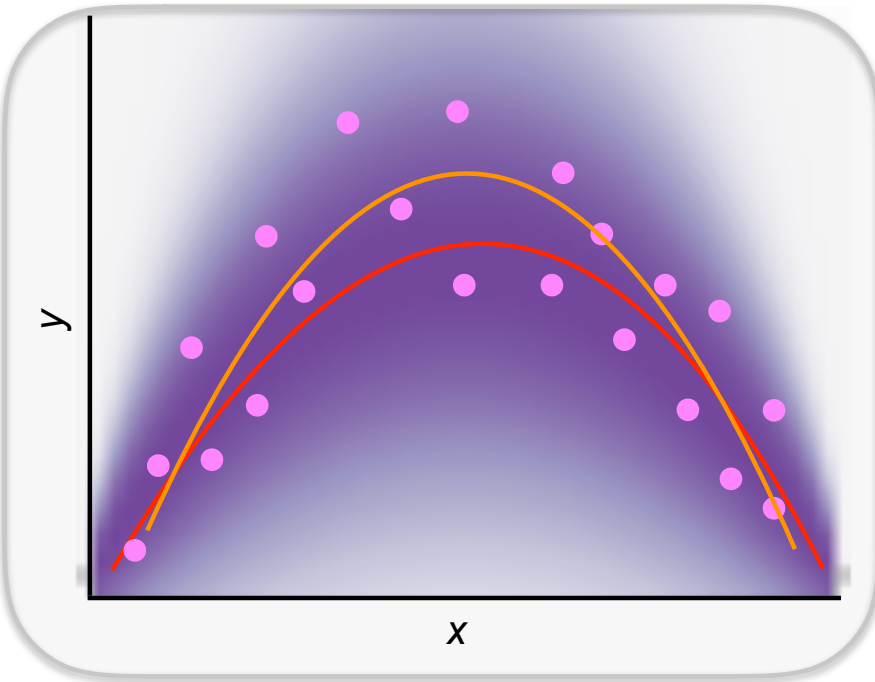
But we only have samples:
 $(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY}$ for $i = 1, \dots, n$

and are restricted to a
function class (e.g., linear)
so we compute:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

Statistical Learning

$$P_{XY}(X = x, Y = y)$$



Ideally, we want to find:

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

But we only have samples:
 $(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY}$ for $i = 1, \dots, n$

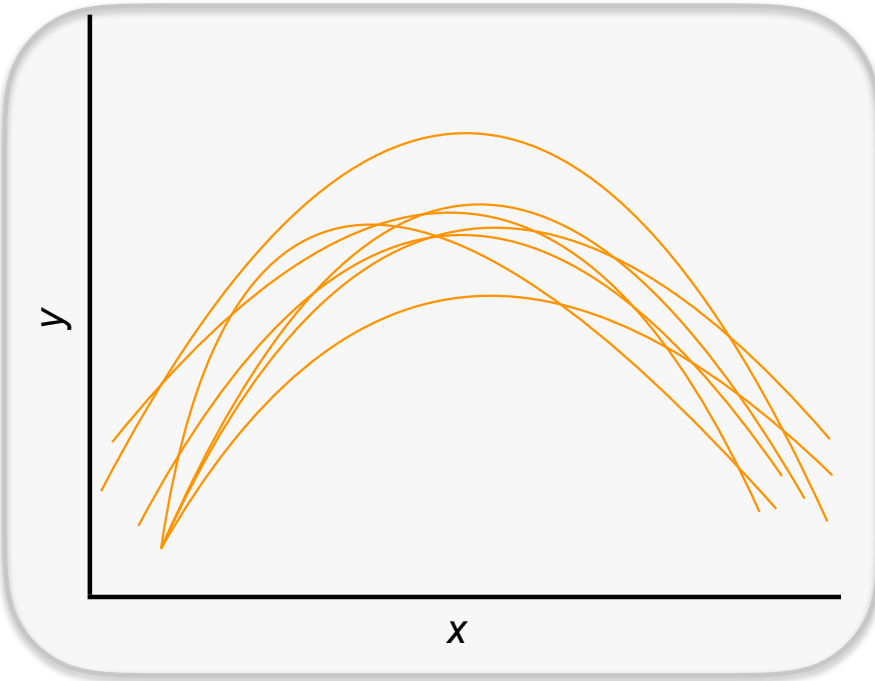
and are restricted to a
function class (e.g., linear)
so we compute:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

We care about future predictions: $\mathbb{E}_{XY}[(Y - \hat{f}(X))^2]$

Statistical Learning

$$P_{XY}(X = x, Y = y)$$



Ideally, we want to find:

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

But we only have samples:

$$(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY} \quad \text{for } i = 1, \dots, n$$

and are restricted to a function class (e.g., linear)

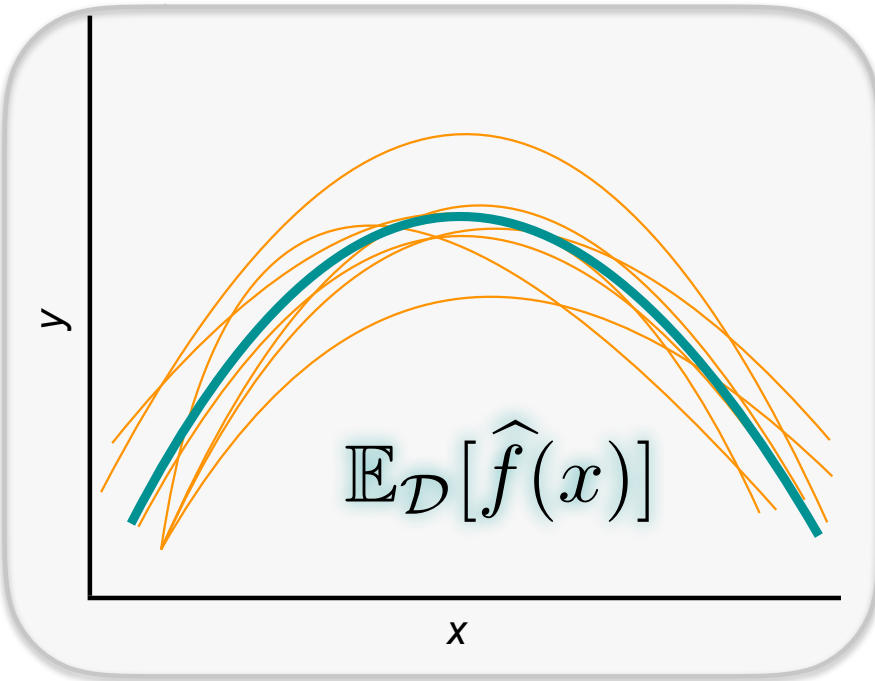
so we compute:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

Each draw $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ results in different \hat{f}

Statistical Learning

$$P_{XY}(X = x, Y = y)$$



Ideally, we want to find:

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

But we only have samples:

$$(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY} \quad \text{for } i = 1, \dots, n$$

and are restricted to a function class (e.g., linear)

so we compute:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

Each draw $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ results in different \hat{f}

Bias-Variance Tradeoff

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x] \qquad \hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

$$\mathbb{E}_{Y|X}[\mathbb{E}_{\mathcal{D}}[(Y - \hat{f}_{\mathcal{D}}(x))^2] | X = x] = \mathbb{E}_{Y|X}[\mathbb{E}_{\mathcal{D}}[(Y - \eta(x) + \eta(x) - \hat{f}_{\mathcal{D}}(x))^2] | X = x]$$

$$= \mathbb{E}_{Y|X} \left[\mathbb{E}_{\mathcal{D}} \left[(Y - \eta(x))^2 + 2(Y - \eta(x))(\eta(x) - \hat{f}_{\mathcal{D}}(x)) + (\eta(x) - \hat{f}_{\mathcal{D}}(x))^2 \middle| X = x \right] \right]$$

$$= \mathbb{E}_{Y|X} \left[(Y - \eta(x))^2 \middle| X = x \right] + 2 \mathbb{E}_{\mathcal{D}} \left[\mathbb{E}_{Y|X} \left[(Y - \eta(x))(\eta(x) - \hat{f}_{\mathcal{D}}(x)) \right] \right] + \mathbb{E}_{\mathcal{D}} \left[(\eta(x) - \hat{f}_{\mathcal{D}}(x))^2 \right]$$

$\mathbb{E}_{Y|X} \left[(Y - \eta(x)) \middle| X = x \right] = \mathbb{E}_{Y|X} \left[Y \middle| X = x \right] - \eta(x) = 0$

Bias-Variance Tradeoff

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x] \qquad \hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

$$\begin{aligned} \mathbb{E}_{Y|X}[\mathbb{E}_{\mathcal{D}}[(Y - \hat{f}_{\mathcal{D}}(x))^2]|X = x] &= \mathbb{E}_{Y|X}[\mathbb{E}_{\mathcal{D}}[(Y - \eta(x) + \eta(x) - \hat{f}_{\mathcal{D}}(x))^2]|X = x] \\ &= \mathbb{E}_{Y|X} \left[\mathbb{E}_{\mathcal{D}}[(Y - \eta(x))^2 + 2(Y - \eta(x))(\eta(x) - \hat{f}_{\mathcal{D}}(x)) \right. \\ &\quad \left. + (\eta(x) - \hat{f}_{\mathcal{D}}(x))^2] | X = x \right] \\ &= \mathbb{E}_{Y|X}[(Y - \eta(x))^2 | X = x] + \mathbb{E}_{\mathcal{D}}[(\eta(x) - \hat{f}_{\mathcal{D}}(x))^2] \end{aligned}$$

irreducible error

Caused by stochastic
label noise

learning error

Caused by either using too “simple”
of a model or not enough
data to learn the model accurately

Bias-Variance Tradeoff

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x] \quad \hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

$$\mathbb{E}_{\mathcal{D}}[(\eta(x) - \hat{f}_{\mathcal{D}}(x))^2] = \mathbb{E}_{\mathcal{D}}[(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] + \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2]$$

$$= \mathbb{E}_{\mathcal{D}}[(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2] + 2 \mathbb{E}_{\mathcal{D}}[(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])(\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))] + \mathbb{E}_{\mathcal{D}}[(\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2]$$

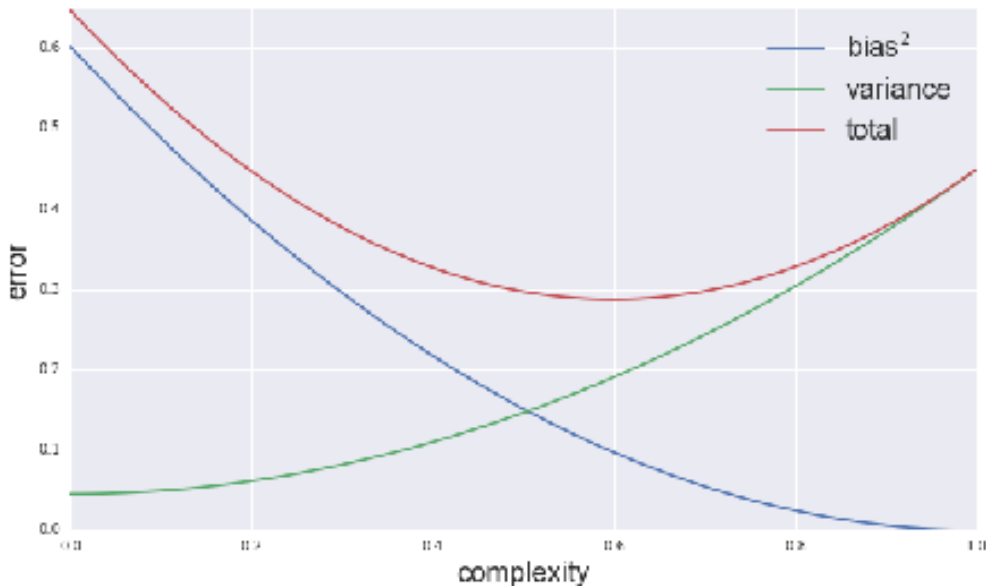
Bias-Variance Tradeoff

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x] \qquad \hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

$$\begin{aligned} \underline{\mathbb{E}_{\mathcal{D}}[(\eta(x) - \hat{f}_{\mathcal{D}}(x))^2]} &= \mathbb{E}_{\mathcal{D}}[(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] + \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2] \\ &= \mathbb{E}_{\mathcal{D}}[(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2 + 2(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])(\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x)) \\ &\quad + (\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2] \\ &= \underline{(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2} + \underline{\mathbb{E}_{\mathcal{D}}[(\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2]} \\ &\qquad \text{biased squared} \qquad \qquad \qquad \text{variance} \end{aligned}$$

Bias-Variance Tradeoff

$$\mathbb{E}_{Y|X}[\mathbb{E}_{\mathcal{D}}[(Y - \hat{f}_{\mathcal{D}}(x))^2] | X = x] = \underbrace{\mathbb{E}_{Y|X}[(Y - \eta(x))^2 | X = x]}_{\text{irreducible error}} + \underbrace{(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2}_{\text{biased squared}} + \underbrace{\mathbb{E}_{\mathcal{D}}[(\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2]}_{\text{variance}}$$



Example: Linear LS

$$\mathbf{Y} = \mathbf{X}w + \epsilon$$

if $y_i = x_i^T w + \epsilon_i$ and $\epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$

$$\hat{w}_{MLE} = \underbrace{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}} = w + \underbrace{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon}$$

$$\underline{\eta(x)} = \mathbb{E}_{Y|X}[Y|X=x] = \mathbb{E}[x^T w + \epsilon | X=x] = \underline{x^T w}$$

$$\hat{f}_D(x) = \hat{w}^T x \quad \underbrace{\mathbb{E}_D[\hat{f}_D(x)]}_{=} = \mathbb{E}_{D|x}[\mathbb{E}_{Y|x}[\hat{w}^T x | X=x]] = \mathbb{E}_{D|x}[w^T x] = w^T x$$

Example: Linear LS $\mathbf{Y} = \mathbf{X}w + \epsilon$

if $y_i = x_i^T w + \epsilon_i$ and $\epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$

$$\hat{w}_{MLE} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon$$

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

$$\hat{f}_{\mathcal{D}}(x) = \hat{w}^T x = w^T x + \epsilon^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} x$$

$$\mathbb{E}_{XY}[\overbrace{(Y - \eta(x))^2}^{\epsilon} | X = x] = \sigma^2$$

irreducible error

$$\underbrace{(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2}_{\text{biased squared}} = 0$$

biased squared

Example: Linear LS $\mathbf{Y} = \mathbf{X}w + \epsilon$

if $y_i = x_i^T w + \epsilon_i$ and $\epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$

$$\hat{w}_{MLE} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon$$

$$\hat{f}_{\mathcal{D}}(x) = \underline{\hat{w}^T x = w^T x + \epsilon^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} x}$$

$$\begin{aligned} \underline{\mathbb{E}_{\mathcal{D}}[(\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2]} &= \mathbb{E}_{\mathcal{D}} \left[x^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \underline{\epsilon \epsilon^T} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} x \right] \\ &= \mathbb{E}_{x_0} \left[x^T (\mathbf{X}^T \mathbf{X})^{-1} x \right] \sigma^2 \end{aligned}$$

variance

Example: Linear LS $\mathbf{Y} = \mathbf{X}w + \epsilon$

$$\text{if } y_i = x_i^T w + \epsilon_i \quad \text{and} \quad \epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$$

$$\hat{w}_{MLE} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon$$

$$\hat{f}_{\mathcal{D}}(x) = \hat{w}^T x = w^T x + \epsilon^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} x$$

$$\begin{aligned} \underbrace{\mathbb{E}_{\mathcal{D}} [(\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2]}_{\text{variance}} &= \mathbb{E}_{\mathcal{D}} [x^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon \epsilon^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} x] \\ &= \sigma^2 x^T (\mathbf{X}^T \mathbf{X})^{-1} x \\ &= \sigma^2 \text{Trace}(\underbrace{(\mathbf{X}^T \mathbf{X})^{-1}} x x^T) \end{aligned}$$

$$\mathbf{X}^T \mathbf{X} = \sum_{i=1}^n x_i x_i^T \stackrel{n \text{ large}}{\rightarrow} n \Sigma \quad \Sigma = \mathbb{E}[X X^T], \quad X \sim P_X$$

$$\mathbb{E}_{X=x} [\mathbb{E}_{\mathcal{D}} [(\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2]] = \frac{\sigma^2}{n} \mathbb{E}_X [\text{Trace}(\Sigma^{-1} X X^T)] = \frac{d\sigma^2}{n}$$

Example: Linear LS $Y = \mathbf{X}w + \epsilon$

if $y_i = x_i^T w + \epsilon_i$ and $\epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$

$$\hat{w}_{MLE} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon$$

$$\eta(x) = \mathbb{E}_{Y|X} [Y | X = x]$$

$$\hat{f}_{\mathcal{D}}(x) = \hat{w}^T x = w^T x + \epsilon^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} x$$

$$\underbrace{\mathbb{E}_{XY} [(Y - \eta(x))^2 | X = x]}_{\text{irreducible error}} = \sigma^2 \qquad \underbrace{(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2}_{\text{biased squared}} = 0$$

$$\mathbb{E}_{X=x} \left[\underbrace{\mathbb{E}_{\mathcal{D}} [(\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2]}_{\text{variance}} \right] = \frac{d\sigma^2}{n}$$