# Linear Regression and the Bias Variance Tradeoff
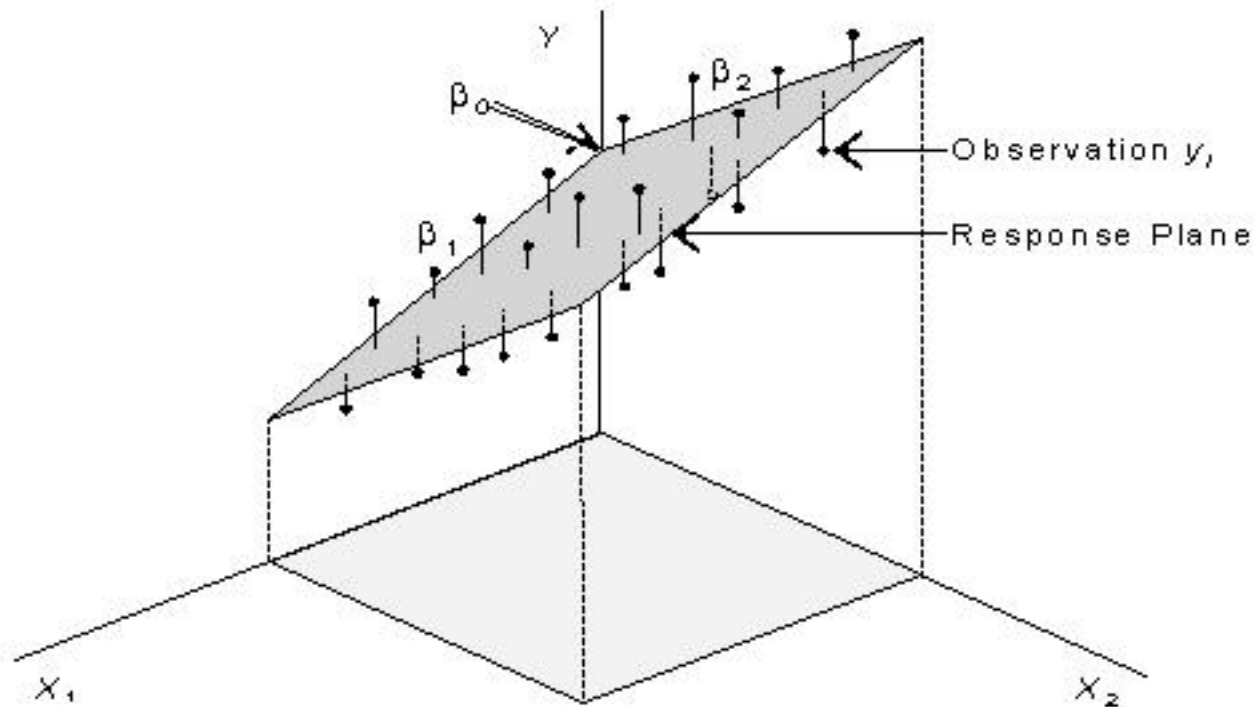
Xiaochun MAI
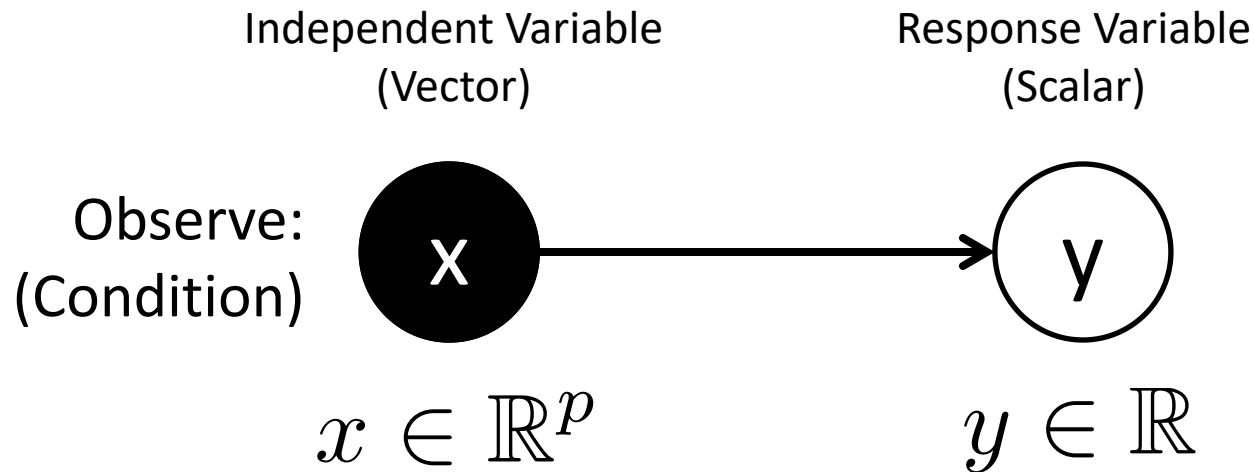
Shenzhen University

# Motivation

- One of the most widely used techniques
- Fundamental to many larger models
  - Generalized Linear Models
  - Collaborative filtering
- Easy to interpret
- Efficient to solve

# Multiple Linear Regression

# The Regression Model

- For a *single* data point *(x,y)*:

Independent Variable
(Vector)

Response Variable
(Scalar)

Observe:
(Condition)

$$x \in \mathbb{R}^p \qquad\qquad y \in \mathbb{R}$$

- Joint Probability:

$$p(x,y) = p(x)\,p(y|x)$$

Discriminative
Model

# The Linear Model

Vector of
Parameters

Vector of
Covariates

Scalar
Response

$$y = \theta^T x + \epsilon$$

Real Value
Noise

$+ \ b$

**Linear Combination**
of Covariates

$$\sum_{i=1}^{p} \theta_i x_i$$

Noise Model:

$$\epsilon \sim N(0, \sigma^2)$$

What about bias/intercept term?

Define: $x_{p+1} = 1$

Then redefine p := p+1 for notational simplicity

# Conditional Likelihood p(y|x)

- Conditioned on x:

$$y = \overbrace{\theta^T x}^{\text{Constant}} + \boxed{\underset{\text{Mean} \quad \text{Variance}}{\epsilon \sim \overset{\text{Normal Distribution}}{N(0, \sigma^2)}}}$$

- Conditional distribution of Y:

$$Y \sim N(\theta^T x, \sigma^2)$$

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

# Parameters and Random Variables
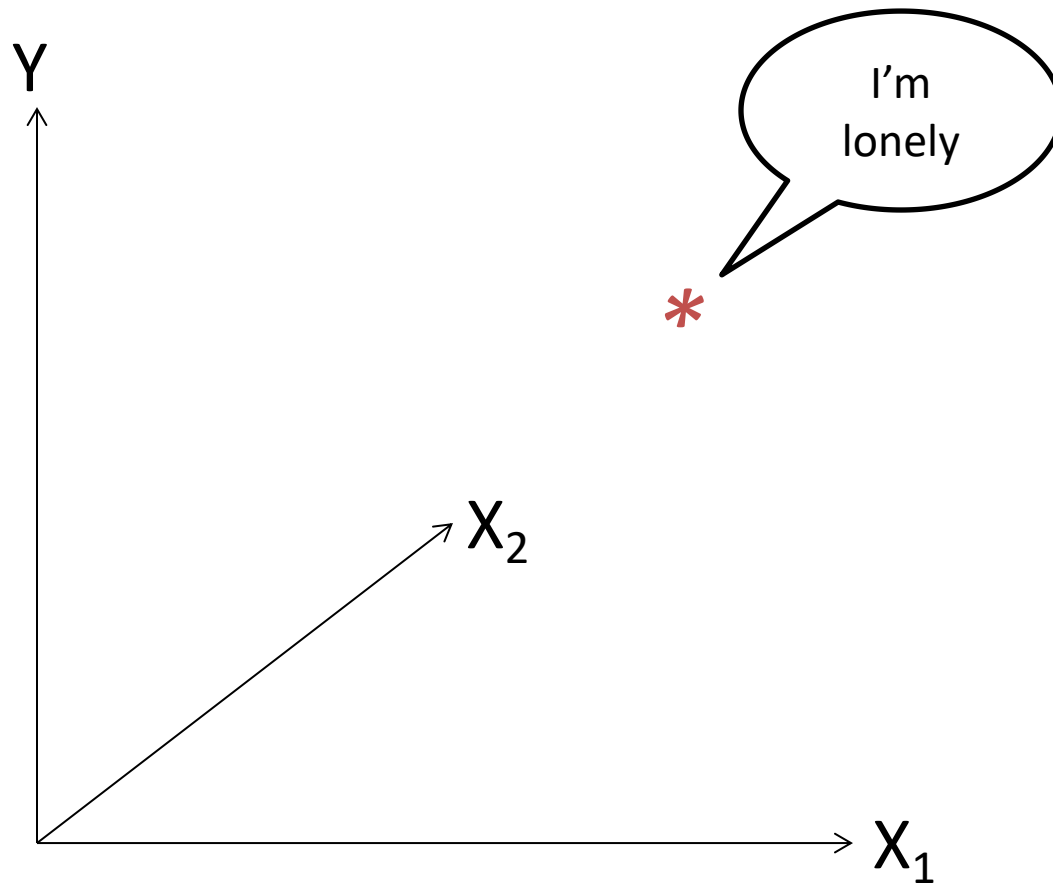
Parameters

$$y \sim N(\theta^T x, \sigma^2)$$

- Conditional distribution of y:
  - Bayesian: parameters as random variables

$$p(y|x, \theta, \sigma^2)$$

  - Frequentist: parameters as (unknown) constants

$$p_{\theta, \sigma^2}(y|x)$$

# So far …

# Independent and Identically Distributed (iid) Data

- For *n* data points:

$$\mathcal{D} = \{(x_1, y_1), \ldots, (x_n, y_n)\}$$
$$= \{(x_i, y_i)\}_{i=1}^{n}$$

Plate Diagram



Independent Variable
(Vector)

Response Variable
(Scalar)

$$x_i \in \mathbb{R}^p \qquad y_i \in \mathbb{R}$$
$$i \in \{1, \ldots, n\}$$

# Joint Probability



- For *n* data points **independent and identically distributed (iid)**:

$$p(\mathcal{D}) = \prod_{i=1}^{n} p(x_i, y_i)$$

$$= \prod_{i=1}^{n} p(x_i)p(y_i|x_i)$$

# Rewriting with Matrix Notation

- Represent data $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{n}$ as:

Covariate (Design) Matrix

Response Vector

$$X = \begin{bmatrix} -\!\!\!- & x_1 & -\!\!\!- \\ -\!\!\!- & x_2 & -\!\!\!- \\ & \cdots & \\ -\!\!\!- & x_n & -\!\!\!- \end{bmatrix} \in \mathbb{R}^{np} \qquad Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^{n}$$

n

p

Assume $X$ has rank p (not degenerate)

n

1

# Rewriting with Matrix Notation

- Rewriting the model using matrix operations:

$$Y = X\theta + \epsilon$$

$$Y \quad = \quad X \quad \theta \quad + \quad \epsilon$$

# Estimating the Model

- Given data how can we estimate θ?

$$Y = X\theta + \epsilon$$

- Construct maximum likelihood estimator (MLE):
  - Derive the log-likelihood
  - Find $\theta_{MLE}$ that maximizes log-likelihood
    - Analytically: Take derivative and set = 0
    - Iteratively: (Stochastic) gradient descent

# Joint Probability



- For *n* data points:

$$p(\mathcal{D}) = \prod_{i=1}^{n} p(x_i, y_i)$$

$$= \prod_{i=1}^{n} \underbrace{p(x_i)}_{\text{"1"}} p(y_i | x_i) \quad \boxed{\text{Discriminative Model}}$$

# Defining the Likelihood



$$p_\theta(y|x) =$$

$$\frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(y - \theta^T x)^2}{2\sigma^2}\right)$$

$$\mathcal{L}(\theta|\mathcal{D}) = \prod_{i=1}^{n} p_\theta(y_i|x_i)$$

$$= \prod_{i=1}^{n} \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(y_i - \theta^T x_i)^2}{2\sigma^2}\right)$$

$$= \frac{1}{\sigma^n (2\pi)^{\frac{n}{2}}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^{n} (y_i - \theta^T x_i)^2\right)$$

# Maximizing the Likelihood

- Want to compute:

$$\hat{\theta}_{\mathrm{MLE}} = \arg \max_{\theta \in \mathbb{R}^p} \mathcal{L}(\theta | \mathcal{D})$$

- To simplify the calculations we take the log:

$$\hat{\theta}_{\mathrm{MLE}} = \arg \max_{\theta \in \mathbb{R}^p} \log \mathcal{L}(\theta | \mathcal{D})$$

which does not affect the maximization because log is a monotone function.

$$\mathcal{L}(\theta|\mathcal{D}) = \frac{1}{\sigma^n (2\pi)^{\frac{n}{2}}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^{n} (y_i - \theta^T x_i)^2\right)$$

- Take the log:

$$\log \mathcal{L}(\theta|\mathcal{D}) = -\log(\sigma^n (2\pi)^{\frac{n}{2}}) - \frac{1}{2\sigma^2} \sum_{i=1}^{n} (y_i - \theta^T x_i)^2$$

- Removing constant terms with respect to θ:

$$\log \mathcal{L}(\theta) = -\sum_{i=1}^{n} (y_i - \theta^T x_i)^2$$

Monotone Function
(Easy to maximize)

$$\log \mathcal{L}(\theta) = -\sum_{i=1}^{n}(y_i - \theta^T x_i)^2$$
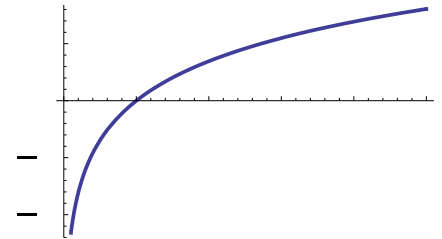
- Want to compute:

$$\hat{\theta}_{\text{MLE}} = \arg\max_{\theta \in \mathbb{R}^p} \log \mathcal{L}(\theta|\mathcal{D})$$

- Plugging in log-likelihood:

$$\hat{\theta}_{\text{MLE}} = \arg\max_{\theta \in \mathbb{R}^p} -\sum_{i=1}^{n}(y_i - \theta^T x_i)^2$$

$$\hat{\theta}_{\mathrm{MLE}} = \arg\max_{\theta \in \mathbb{R}^p} -\sum_{i=1}^{n}(y_i - \theta^T x_i)^2$$

- Dropping the sign and flipping from maximization to minimization:

$$\hat{\theta}_{\mathrm{MLE}} = \arg\min_{\theta \in \mathbb{R}^p} \underbrace{\sum_{i=1}^{n}(y_i - \theta^T x_i)^2}$$

Minimize Sum (Error)$^2$

- Gaussian Noise Model → Squared Loss
  - Least Squares Regression

# Pictorial Interpretation of Squared Error

# Maximizing the Likelihood (Minimizing the Squared Error)

$$\hat{\theta}_{\mathrm{MLE}} = \arg \min_{\theta \in \mathbb{R}^p} \sum_{i=1}^{n} (y_i - \theta^T x_i)^2$$

Convex Function

$-\log \mathcal{L}(\theta)$

Slope = 0

$\theta$

$\hat{\theta}_{\mathrm{MLE}}$

- Take the gradient and set it equal to zero

# Minimizing the Squared Error

$$\hat{\theta}_{\mathrm{MLE}} = \arg\min_{\theta \in \mathbb{R}^p} \sum_{i=1}^{n} (y_i - \theta^T x_i)^2$$

- Taking the gradient

$$-\nabla_\theta \log \mathcal{L}(\theta) = \nabla_\theta \sum_{i=1}^{n} (y_i - \theta^T x_i)^2$$

Chain Rule →
$$= -2 \sum_{i=1}^{n} (y_i - \theta^T x_i) x_i$$

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

- Rewriting the gradient in matrix form:

$$-\nabla_\theta \log \mathcal{L}(\theta) = -2 \sum_{i=1}^{n} y_i x_i + 2 \sum_{i=1}^{n} (\theta^T x_i) x_i$$

$$= -2X^T Y + 2X^T X \theta$$

- To make sure the log-likelihood is convex compute the second derivative (Hessian)

$$-\nabla^2 \log \mathcal{L}(\theta) = 2X^T X$$

- If $X$ is full rank then $X^T X$ is positive definite and therefore $\theta_{MLE}$ is the minimum.
  - Address the degenerate cases with regularization

$$-\nabla_\theta \log \mathcal{L}(\theta) = -2X^T y + 2X^T X \theta = 0$$

- Setting gradient equal to 0 and solve for $\theta_{\text{MLE}}$:

$$(X^T X)\hat{\theta}_{\text{MLE}} = X^T Y$$

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

Normal Equations
(Write on board)

# Geometric Interpretation

- View the MLE as finding a projection on col(X)
  - Define the estimator:
    $$\hat{Y} = X\theta$$
  - Observe that Ŷ is in col(X)
    - linear combination of cols of X
  - Want to Ŷ closest to Y
- Implies (Y-Ŷ) normal to X

$$X^T(Y - \hat{Y}) = X^T(Y - X\theta) = 0$$
$$\Rightarrow X^T X \theta = X^T Y$$

# Connection to Pseudo-Inverse

$$\hat{\theta}_{\mathrm{MLE}} = \underbrace{(X^T X)^{-1} X^T}_{} Y$$

Moore-Penrose $X^\dagger$
Pseudoinverse

- Generalization of the inverse:
  - Consider the case when X is square and invertible:

$$X^\dagger = (X^T X)^{-1} X^T = X^{-1} (X^T)^{-1} X^T = X^{-1}$$

  - Which implies $\theta_{\mathrm{MLE}} = X^{-1} Y$ the solution to $X\theta = Y$ when $X$ is square and invertible

# Computing the MLE

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

- **Not** typically solved by inverting $X^T X$
- Solved using direct methods:
  - Cholesky factorization:
    - Up to a factor of 2 faster
  - QR factorization:
    - More numerically stable

  or use the
  built-in solver
  in your math library.
  R: solve(Xt %*% X, Xt %*% y)

- Solved using various iterative methods:
  - Krylov subspace methods
  - (Stochastic) Gradient Descent

# Cholesky Factorization

$$\text{solve} \quad (\underbrace{X^T X}_{C})\hat{\theta}_{\text{MLE}} = \underbrace{X^T Y}_{d}$$
$$\hat{\theta}_{\text{MLE}}$$

- Compute symm. matrix $C = X^T X$ $\qquad O(np^2)$
- Compute vector $d = X^T Y$ $\qquad O(np)$
- Cholesky Factorization $LL^T = C$ $\qquad O(p^3)$
  - L is lower triangular
- Forward subs. to solve: $Lz = d$ $\qquad O(p^2)$
- Backward subs. to solve: $L^T \hat{\theta}_{\text{MLE}} = z$ $\qquad O(p^2)$

Connections to graphical model inference:
http://ssg.mit.edu/~willsky/publ_pdfs/185_pub_MLR.pdf and http://yaroslavvb.blogspot.com/2011/02/junction-trees-in-numerical-analysis.html with illustrations

# Solving Triangular System

$$
\begin{array}{|c|c|c|c|}
\hline
A_{11}x_1 & A_{12}x_2 & A_{13}x_3 & A_{14}x_4 \\
\hline
 & A_{22} & A_{23} & A_{24} \\
\hline
 & & A_{33} & A_{34} \\
\hline
 & & & A_{44} \\
\hline
\end{array}
\ast
\begin{array}{|c|}
\hline
x_1 \\
\hline
x_2 \\
\hline
x_3 \\
\hline
x_4 \\
\hline
\end{array}
=
\begin{array}{|c|}
\hline
b_1 \\
\hline
b_2 \\
\hline
b_3 \\
\hline
b_4 \\
\hline
\end{array}
$$

Bonus Content

# Solving Triangular System

| $A_{11}x_1$ | $A_{12}x_2$ | $A_{13}x_3$ | $A_{14}x_4$ | | $b_1$ |
|---|---|---|---|---|---|
| | $A_{22}x_2$ | $A_{23}x_3$ | $A_{24}x_4$ | | $b_2$ |
| | | $A_{33}x_3$ | $A_{34}x_4$ | | $b_3$ |
| | | | $A_{44}x_4$ | | $b_4$ |

$$x_1 = \frac{b_1 - A_{12}x_2 - A_{13}x_3 - A_{14}x_4}{A_{11}}$$

$$x_2 = \frac{b_2 - A_{23}x_3 - A_{24}x_4}{A_{22}}$$

$$x_3 = \frac{(b_3 - A_{34}x_4)}{A_{33}}$$

$$x_4 = b_4 / A_{44}$$

Bonus Content

# Distributed Direct Solution (Map-Reduce)

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

- Distribution computations of sums:

p

p $\quad C = X^T X = \sum_{i=1}^{n} x_i x_i^T \qquad O(np^2)$

1

p $\quad d = X^T y = \sum_{i=1}^{n} x_i y_i \qquad O(np)$

- Solve system $C\,\theta_{\mathrm{MLE}} = d$ on master. $\qquad O(p^3)$

# Gradient Descent:
## What if p is large?  (e.g., n/2)

- The cost of O($np^2$) = O($n^3$) could by prohibitive

- Solution: Iterative Methods
  - Gradient Descent:

For $\tau$ from 0 until *convergence*

$$\theta^{(\tau+1)} = \theta^{(\tau)} - \rho(\tau)\nabla \log \mathcal{L}(\theta^{(\tau)}|D)$$

Learning rate

# Gradient Descent Illustrated:

$-\log \mathcal{L}(\theta)$

$\theta^{(0)}$

Slope = 0

$\theta^{(2)}$ $\theta^{(3)}$ $\theta^{(1)}$

Convex Function

$\theta^{(3)} = \hat{\theta}_{\mathrm{MLE}}$

$\theta$

# Gradient Descent:
## What if p is large?  (e.g., n/2)

- The cost of O($np^2$) = O($n^3$) could by prohibitive
- Solution: Iterative Methods
  - Gradient Descent:

For $\tau$ from 0 until *convergence*

$$\theta^{(\tau+1)} = \theta^{(\tau)} - \rho(\tau)(-\nabla \log \mathcal{L}(\theta^{(\tau)}|D))$$

$$= \theta^{(\tau)} + \rho(\tau)\frac{1}{n}\sum_{i=1}^{n}(y_i - \theta^{(\tau)T}x_i)x_i \quad O(np)$$

- Can we do better?

Estimate of the Gradient

# Supplement: Derivation Process

$$\hat{\theta}_{\mathrm{MLE}} = \arg \min_{\theta \in \mathbb{R}^p} \sum_{i=1}^{n} (y_i - \theta^T x_i)^2$$

$$= arg\,min_{\theta \in \mathbb{R}^p} \quad \frac{1}{2n} \sum_{i=1}^{n} (y_i - \theta^T x_i)^2$$

Larger $x_i$, larger the gradients. To avoid the increasing of the gradients with the number of $x_i$, multiply a 1/2n.

$$-\nabla_\theta \log(\theta) = \frac{1}{2n} \sum_{i=1}^{n} 2(y_i - \theta^T x_i)(-x_i)$$

$$= -\frac{1}{n} \sum_{i=1}^{n} (y_i - \theta^T x_i) x_i$$

# Stochastic Gradient Descent

- Construct noisy estimate of the gradient:

For $\tau$ from 0 until *convergence*
   *1) pick a random i*
   *2)* $\theta^{(\tau+1)} = \theta^{(\tau)} + \rho(\tau)\boxed{(y_i - \theta^{(\tau)T}x_i)x_i}$    $O(p)$

- Sensitive to choice of ρ(τ) typically (ρ(τ)=1/τ)

- Also known as Least-Mean-Squares (LMS)

- Applies to streaming data *O(p)* storage

# Fitting Non-linear Data

- What if Y has a non-linear response?



- Can we still use a linear model?

# Transforming the Feature Space

- Transform features *x<sub>i</sub>*

$$x_i = (x_{i,1}, x_{i,2}, \ldots, x_{i,p})$$

- By applying non-linear transformation $\phi$:

$$\phi : \mathbb{R}^p \to \mathbb{R}^k$$

- Example:

$$\phi(x) = \{1, x, x^2, \ldots, x^k\}$$

  - others: splines, radial basis functions, …
  - Expert engineered features (modeling)

# Under-fitting vs over-fitting



https://scikit-learn.org/0.15/auto_examples/plot_underfitting_overfitting.html
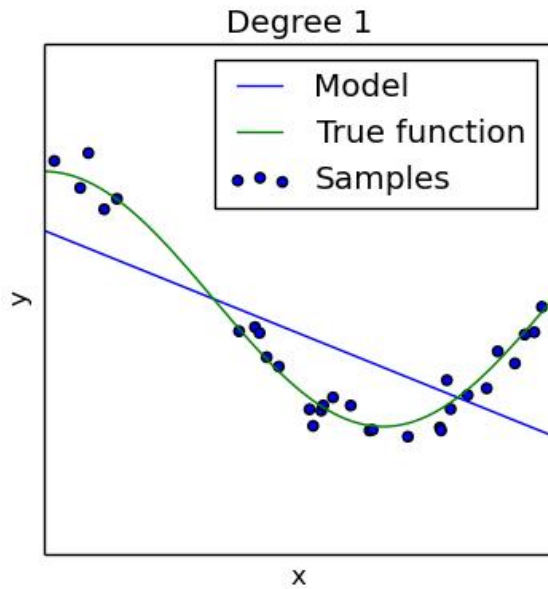
# Really Over-fitting!



- Errors on training data are small
- But errors on new points are likely to be large

# Bias-Variance Tradeoff

- So far we have minimized the error (loss) with respect to **training data**

  - Low training error does not imply good expected performance: **over-fitting**

- We would like to reason about the **expected loss (Prediction Risk)** over:

  - Training Data: $\{(y_1, x_1), ..., (y_n, x_n)\}$
  - Test point: $(y_*, x_*)$

- We will decompose the expected loss into:

$$\mathbf{E}_{D,(y_*,x_*)}\left[(y_* - f(x_*|D))^2\right] = \mathrm{Noise} + \mathrm{Bias}^2 + \mathrm{Variance}$$

# What if I train on different data?

Low Variability:



High Variability

- Define (unobserved) the true model (*h*):

$$y_* = h(x_*) + \epsilon_*$$

Assume 0 mean noise
[bias goes in *h(x*)*]

- Completed the squares with: $h(x_*) = h_*$

$$\mathbf{E}_{D,(y_*,x_*)} \left[ (y_* - f(x_*|D))^2 \right]$$ Expected Loss

$$= \mathbf{E}_{D,(y_*,x_*)} \left[ (y_* \underbrace{- h(x_*)}_{a} + \underbrace{h(x_*) - f(x_*|D)}_{b})^2 \right]$$

$$(a+b)^2 = a^2 + b^2 + 2ab$$

$$= \mathbf{E}_{\epsilon_*} \left[ (y_* - h(x_*))^2 \right] + \mathbf{E}_D \left[ (h(x_*) - f(x_*|D))^2 \right]$$
$$+ 2\mathbf{E}_{D,(y_*,x_*)} \left[ y_* h_* - y_* f_* - h_* h_* + h_* f_* \right]$$

- Define (unobserved) the true model (*h*):

$$y_* = h(x_*) + \epsilon_*$$

- Completed the squares with: $h(x_*) = h_*$

$$\mathbf{E}_{D,(y_*,x_*)}\left[(y_* - f(x_*|D))^2\right] \quad \text{Expected Loss}$$

$$= \mathbf{E}_{D,(y_*,x_*)}\left[(y_* - h(x_*) + h(x_*) - f(x_*|D))^2\right]$$

$$= \mathbf{E}_{\epsilon_*}\left[(y_* - h(x_*))^2\right] + \mathbf{E}_D\left[(h(x_*) - f(x_*|D))^2\right]$$

$$+ 2\mathbf{E}_{D,(y_*,x_*)}\left[y_* h_* - y_* f_* - h_* h_* + h_* f_*\right]$$

Substitute defn. $y_* = h_* + e_*$

$$\mathbf{E}\left[(h_* + \epsilon_*)h_* - (h_* + \epsilon_*)f_* - h_* h_* + h_* f_*\right] =$$

$$h_* h_* + \mathbf{E}[\epsilon_*]h_* - h_*\mathbf{E}[f_*] - \mathbf{E}[\epsilon_*]f_* - h_* h_* + h_*\mathbf{E}[f_*]$$

- Define (unobserved) the true model (*h*):

$$y_* = h(x_*) + \epsilon_*$$

- Completed the squares with: $h(x_*) = h_*$

$$\mathbf{E}_{D,(y_*,x_*)}\left[(y_* - f(x_*|D))^2\right] \text{ Expected Loss}$$

$$= \mathbf{E}_{D,(y_*,x_*)}\left[(y_* - h(x_*) + h(x_*) - f(x_*|D))^2\right]$$

$$= \mathbf{E}_{\epsilon_*}\left[(y_* - h(x_*))^2\right] + \mathbf{E}_D\left[(h(x_*) - f(x_*|D))^2\right]$$

Noise Term
(out of our control)
☹

Model Estimation Error
(we want to minimize this)

Expand

- Minimum error is governed by the noise.

- Expanding on the model estimation error:

$$\mathbf{E}_D\left[(h(x_*) - f(x_*|D))^2\right]$$

- Completing the squares with $\mathbf{E}\left[f(x_*|D)\right] = \bar{f}_*$

$$\mathbf{E}_D\left[(h(x_*) - f(x_*|D))^2\right]$$

$$= \mathbf{E}\left[(h(x_*) - \mathbf{E}\left[f(x_*|D)\right] + \mathbf{E}\left[f(x_*|D)\right] - f(x_*|D))^2\right]$$

$$= \mathbf{E}\left[(h(x_*) - \mathbf{E}\left[f(x_*|D)\right])^2\right] + \mathbf{E}\left[(f(x_*|D) - \mathbf{E}\left[f(x_*|D)\right])^2\right]$$

$$+ 2\mathbf{E}\left[h_*\bar{f}_* - h_*f_* - \bar{f}_*f_* + \bar{f}_*^2\right]$$

$$= h_*\bar{f}_* - h_*\mathbf{E}\left[f_*\right] - \bar{f}_*\mathbf{E}\left[f_*\right] + \bar{f}_*^2 =$$

$$h_*\bar{f}_* - h_*\bar{f}_* - \bar{f}_*\bar{f}_* + \bar{f}_*^2 = 0$$

- Expanding on the model estimation error:
$$\mathbf{E}_D \left[ (h(x_*) - f(x_*|D))^2 \right]$$

- Completing the squares with $\mathbf{E}\left[f(x_*|D)\right] = \bar{f}_*$

$$\mathbf{E}_D \left[ (h(x_*) - f(x_*|D))^2 \right]$$
$$= \mathbf{E}\left[ (h(x_*) - \mathbf{E}\left[f(x_*|D)\right])^2 \right] + \mathbf{E}\left[ (f(x_*|D) - \mathbf{E}\left[f(x_*|D)\right])^2 \right]$$

$$(h(x_*) - \mathbf{E}\left[f(x_*|D)\right])^2$$

- Expanding on the model estimation error:

$$\mathbf{E}_D\left[(h(x_*) - f(x_*|D))^2\right]$$

- Completing the squares with $\mathbf{E}\left[f(x_*|D)\right] = \bar{f}_*$

$$\mathbf{E}_D\left[(h(x_*) - f(x_*|D))^2\right]$$
$$= (h(x_*) - \mathbf{E}\left[f(x_*|D)\right])^2 + \mathbf{E}\left[(f(x_*|D) - \mathbf{E}\left[f(x_*|D)\right])^2\right]$$

$(\text{Bias})^2$ \qquad\qquad Variance

- Tradeoff between bias and variance:
  - **Simple Models:** High Bias, Low Variance
  - **Complex Models:** Low Bias, High Variance

# Summary of Bias Variance Tradeoff

$$\mathbf{E}_{D,(y_*,x_*)} \left[ (y_* - f(x_*|D))^2 \right] = \qquad \text{Expected Loss}$$

$$\mathbf{E}_{\epsilon_*} \left[ (y_* - h(x_*))^2 \right] \qquad \text{Noise}$$

$$+ (h(x_*) - \mathbf{E}_D \left[ f(x_*|D) \right])^2 \qquad \text{(Bias)}^2$$

$$+ \mathbf{E}_D \left[ (f(x_*|D) - \mathbf{E}_D \left[ f(x_*|D) \right])^2 \right] \quad \text{Variance}$$

- Choice of models balances bias and variance.
  - Over-fitting ➜ Variance is too High
  - Under-fitting ➜ Bias is too High

# Bias Variance Plot



Image from http://scott.fortmann-roe.com/docs/BiasVariance.html

# Analyze bias of $f(x_*|D) = x_*^T \hat{\theta}_{\mathrm{MLE}}$

- Assume a true model is linear: $h(x_*) = x_*^T \theta$

$$\mathrm{bias} = h(x_*) - \mathbf{E}_D\left[f(x_*|D)\right]$$

$$= x_*^T \theta - \mathbf{E}_D\left[x_*^T \hat{\theta}_{\mathrm{MLE}}\right]$$

$$= x_*^T \theta - \mathbf{E}_D\left[x_*^T (X^T X)^{-1} X^T Y\right]$$

$$= x_*^T \theta - \mathbf{E}_D\left[x_*^T (X^T X)^{-1} X^T (X\theta + \epsilon)\right]$$

$$= x_*^T \theta - \mathbf{E}_D\left[x_*^T (X^T X)^{-1} X^T X\theta + x_*^T (X^T X)^{-1} X^T \epsilon\right]$$

$$= x_*^T \theta - \mathbf{E}_D\left[x_*^T \theta + x_*^T (X^T X)^{-1} X^T \epsilon\right]$$

$$= x_*^T \theta - x_*^T \theta + x_*^T (X^T X)^{-1} X^T \mathbf{E}_D\left[\epsilon\right]$$

$$= x_*^T \theta - x_*^T \theta = 0$$

Substitute MLE

Plug in definition of Y

Expand and cancel

Assumption:
$$\mathbf{E}_D\left[\epsilon\right] = 0$$

$\hat{\theta}_{\mathrm{MLE}}$ is unbiased!

# Analyze Variance of $f(x_*|D) = x_*^T \hat{\theta}_{\mathrm{MLE}}$

- Assume a true model is linear: $h(x_*) = x_*^T \theta$

$$\mathrm{Var.} = \mathbf{E}\left[(f(x_*|D) - \mathbf{E}_D\left[f(x_*|D)\right])^2\right]$$

$$= \mathbf{E}\left[(x_*^T \hat{\theta}_{\mathrm{MLE}} - x_*^T \theta)^2\right]$$

Substitute MLE + unbiased result

$$= \mathbf{E}\left[(x_*^T (X^T X)^{-1} X^T Y - x_*^T \theta)^2\right]$$

Plug in definition of Y

$$= \mathbf{E}\left[(x_*^T (X^T X)^{-1} X^T (X\theta + \epsilon) - x_*^T \theta)^2\right]$$

$$= \mathbf{E}\left[(x_*^T \theta + x_*^T (X^T X)^{-1} X^T \epsilon - x_*^T \theta)^2\right]$$

$$= \mathbf{E}\left[(x_*^T (X^T X)^{-1} X^T \epsilon)^2\right]$$

Expand and cancel

- Use property of scalar: $a^2 = a\, a^T$

# Analyze Variance of $f(x_*|D) = x_*^T \hat{\theta}_{\mathrm{MLE}}$

- Use property of scalar: a² = a aᵀ

$$\mathrm{Var.} = \mathbf{E}\left[(f(x_*|D) - \mathbf{E}_D\left[f(x_*|D)\right])^2\right]$$

$$= \mathbf{E}\left[(x_*^T(X^TX)^{-1}X^T\epsilon)^2\right]$$

$$= \mathbf{E}\left[(x_*^T(X^TX)^{-1}X^T\epsilon)(x_*^T(X^TX)^{-1}X^T\epsilon)^T\right]$$

$$= \mathbf{E}\left[x_*^T(X^TX)^{-1}X^T\epsilon\epsilon^T(x_*^T(X^TX)^{-1}X^T)^T\right]$$

$$= x_*^T(X^TX)^{-1}X^T\mathbf{E}\left[\epsilon\epsilon^T\right](x_*^T(X^TX)^{-1}X^T)^T$$

$$= x_*^T(X^TX)^{-1}X^T\sigma_\epsilon^2 I(x_*^T(X^TX)^{-1}X^T)^T$$

$$= \sigma_\epsilon^2 x_*^T(X^TX)^{-1}X^TX(x_*^T(X^TX)^{-1})^T$$

$$= \sigma_\epsilon^2 x_*^T(x_*^T(X^TX)^{-1})^T$$

$$= \sigma_\epsilon^2 x_*^T(X^TX)^{-1}x_*$$

# Analyze Variance of $f(x_*|D) = x_*^T \hat{\theta}_{\mathrm{MLE}}$

- Assume a true model is linear: $h(x_*) = x_*^T \theta$

$$\mathrm{Var.} = \mathbf{Var}\left[f(x_*|D) - \mathbf{E}_D\left[f(x_*|D)\right]\right]$$

$$= \mathbf{Var}\left[x_*^T \hat{\theta}_{\mathrm{MLE}} - x_*^T \theta\right]$$

Substitute MLE + unbiased result

$$= \mathbf{Var}\left[x_*^T (X^T X)^{-1} X^T Y - x_*^T \theta\right]$$

Plug in definition of Y

$$= \mathbf{Var}\left[x_*^T (X^T X)^{-1} X^T (X\theta + \epsilon) - x_*^T \theta\right]$$

$$= \mathbf{Var}\left[x_*^T \theta + x_*^T (X^T X)^{-1} X^T \epsilon - x_*^T \theta\right]$$

$$= \mathbf{Var}\left[x_*^T (X^T X)^{-1} X^T \epsilon\right]$$

Expand and cancel

- Next: use matrix variance identity

# Analyze Variance of $f(x_* | D) = x_*^T \hat{\theta}_{\text{MLE}}$

- Define: $A = x_*^T (X^T X)^{-1} X^T$

  $$\text{Var.} = \mathbf{Var}\left[ x_*^T (X^T X)^{-1} X^T \epsilon \right] = \mathbf{Var}\left[ A\epsilon \right]$$

- Use matrix variance identity: $\mathbf{Var}\left[ A\epsilon \right] = A\Sigma_\epsilon A^T$

  $$\text{Var.} = A\Sigma_\epsilon A^T = \sigma_\epsilon^2 A A^T$$

  $$= \sigma_\epsilon^2 x_*^T (X^T X)^{-1} X^T (x_*^T (X^T X)^{-1} X^T)^T$$

  $$= \sigma_\epsilon^2 x_*^T (X^T X)^{-1} X^T X (x_*^T (X^T X)^{-1})^T$$

  $$= \sigma_\epsilon^2 x_*^T (x_*^T (X^T X)^{-1})^T$$

  $$= \sigma_\epsilon^2 x_*^T (X^T X)^{-1} x_*$$

- If we assume $x$ is iid N(0, 1): $\mathbf{E}_{X, x_*}\left[ \text{Var.} \right] = \sigma_\epsilon^2 \frac{p}{n}$

# Deriving the final identity

- Assume $x_i$ and $x_*$ are N(0,1)

$$\mathbf{E}_{X,x_*}\left[\text{Var.}\right] = \sigma_\epsilon^2 \mathbf{E}_{X,x_*}\left[x_*^T (X^T X)^{-1} x_*\right]$$
$$= \sigma_\epsilon^2 \mathbf{E}_{X,x_*}\left[tr(x_* x_*^T (X^T X)^{-1})\right]$$
$$= \sigma_\epsilon^2 tr(\mathbf{E}_{X,x_*}\left[x_* x_*^T (X^T X)^{-1}\right])$$
$$= \sigma_\epsilon^2 tr(\mathbf{E}_{x_*}\left[x_* x_*^T\right] \mathbf{E}_X\left[(X^T X)^{-1}\right])$$
$$= \frac{\sigma_\epsilon^2}{n} tr(\mathbf{E}_{x_*}\left[x_* x_*^T\right])$$
$$= \frac{\sigma_\epsilon^2}{n} p$$

# Summary

- Least-Square Regression is Unbiased:

$$\mathbf{E}_D \left[ x_*^T \hat{\theta}_{\mathrm{MLE}} \right] = x_*^T \theta$$

- Variance depends on:

$$\mathbf{E} \left[ (f(x_*|D) - \mathbf{E} \left[ f(x_*|D) \right])^2 \right] = \sigma_\epsilon^2 x_*^T (X^T X)^{-1} x_*$$

$$\approx \sigma_\epsilon^2 \frac{p}{n}$$

- Number of data-points *n*
- Dimensionality *p*
- Not on observations *Y*

# Gauss-Markov Theorem

- The linear model:

$$f(x_*) = x_*^T \hat{\theta}_{\mathrm{MLE}} = x_*^T (X^T X)^{-1} X^T Y$$

has the **minimum variance** among all **unbiased** linear estimators

  – Note that this is linear in Y

- **BLUE: B**est **L**inear **U**nbiased **E**stimator

# Summary

- Introduced the Least-Square regression model
  - Maximum Likelihood: Gaussian Noise
  - Loss Function: Squared Error
  - Geometric Interpretation: Minimizing Projection
- Derived the normal equations:
  - Walked through process of constructing MLE
  - Discussed efficient computation of the MLE
- Introduced basis functions for non-linearity
  - Demonstrated issues with over-fitting
- Derived the classic bias-variance tradeoff
  - Applied to least-squares model

# Additional Reading I found Helpful

- http://www.stat.cmu.edu/~roeder/stat707/lectures.pdf

- http://people.stern.nyu.edu/wgreene/MathStat/GreeneChapter4.pdf

- http://www.seas.ucla.edu/~vandenbe/103/lectures/qr.pdf

- http://www.cs.berkeley.edu/~jduchi/projects/matrix_prop.pdf