# Linear Regression and Bias-Variance Tradeoff

Machine Learning - 10601

Geoff Gordon, Miroslav Dudík
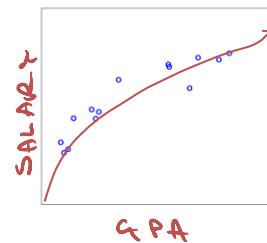
http://www.cs.cmu.edu/~ggordon/10601/
September 21, 2009

---

# Last time: linear regression

**Goal:** TASK

predict a <u>continuous</u> response
from <u>continuous/categorical inputs</u>

EXPERIENCE    (GPA, GENDER ...)

**Input:**
$(\mathbf{x_1}, \mathbf{y_1}), \dots, (\mathbf{x_N}, \mathbf{y_N})$ — SALARY "RESPONSE"

**Model:** LINEAR
$$y \approx \sum_{j=1}^{M} w_j \phi_j(x) \leftarrow \text{FEATURES / BASIS FCTS}$$

**Performance measure:**
$$\sum_n \left( y_n - w \cdot \phi(x_n) \right)^2 \leftarrow \text{MEAN SQUARED ERROR}$$



SALARY vs GPA

## Linear regression

Input: $(\mathbf{x_1}, \mathbf{y_1}), \ldots, (\mathbf{x_N}, \mathbf{y_N})$

Assume: $\mathbf{y} \approx \sum_j \mathbf{w_j \phi_j(x)}$

Goal:
$$\min_w \sum_n \left( y_n - w \cdot \phi(x_n) \right)^2$$

feature vector
inner product

$$\begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix} \approx \begin{pmatrix} \phi(x_1) \cdot w \\ \vdots \\ \phi(x_N) \cdot w \end{pmatrix} = \begin{pmatrix} \phi(x_1) \\ \vdots \\ \phi(x_N) \end{pmatrix} \cdot w$$

ROW VECTOR OF FEATURES

$\Phi$ MATRIX

COLUMN VECTOR OF WEIGHTS

---

LIN. IN $w$

$$\min_w \| y - \underline{\Phi} w \|^2$$

CONVEX IN $w$

VEC. of predictions

$$\begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix}$$

looking for $w$
$\hat{y}$ in column space of $\underline{\Phi}$ that is closest to $y$

$$\nabla_w = -2 \Phi^T (y - \underline{\Phi} w) = 0$$

$$2 \underline{\Phi}^T \underline{\Phi} w = 2 \Phi^T y$$
IF INVERTIBLE

$$w = (\Phi^T \underline{\Phi})^{-1} \Phi^T y$$

PROBLEMS
I. COLINEARITY AMONG FEATURES
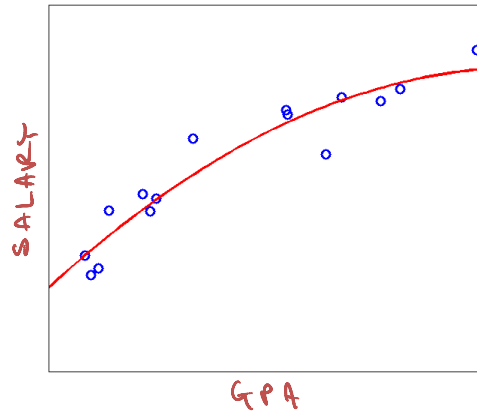II. TOO MANY FEATURES CAN FIT ANYTHING

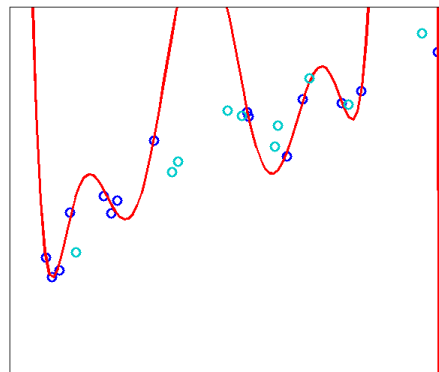INVERSION NOT POSSIBLE IF FEATURES LINEARLY DEPENDENT

# Linear regression
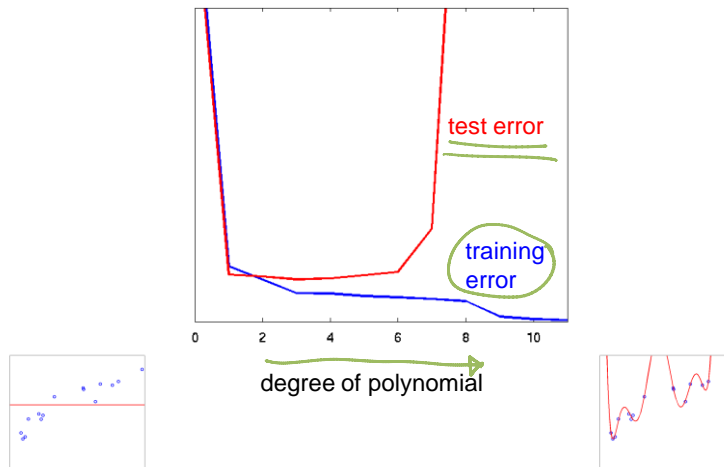
$$y \approx w_0 + w_1 x + w_2 x^2$$



# Linear regression

$$y \approx w_0 + w_1 x + w_2 x^2$$
$$+ w_3 x^3 + w_4 x^4$$
$$+ \dots + w_{10} x^{10}$$



BIG ERROR ON TEST
EVEN THOUGH ZERO
ON TRAINING

# Training & Test Error
# Don't Match



test error

training error

degree of polynomial

# Training & Test Error
# Don't Match: Why?

$$\text{error}_{\text{train}} = \frac{1}{N} \sum_n \left( y_n - \phi(x_n) \cdot w \right)^2$$

unknown dist. on $(x, y)$

ALWAYS TOO OPTIMISTIC

$$\text{error}_{\text{true}} = E_{x, y} \left[ (y - \phi(x) \cdot w)^2 \right]$$

samples $(x_1', y_1'), \dots, (x_K', y_K')$

$$= \int (y - \phi(x) \cdot w)^2 \, p(x, y) \, dx \, dy$$

MONTE CARLO SAMPLING

$$\approx \frac{1}{K} \sum_k \left( y_k' - \phi(x_k') \cdot w \right)^2 = \text{error test}$$

indep. of $\hat{w}$

# Training & Test Error
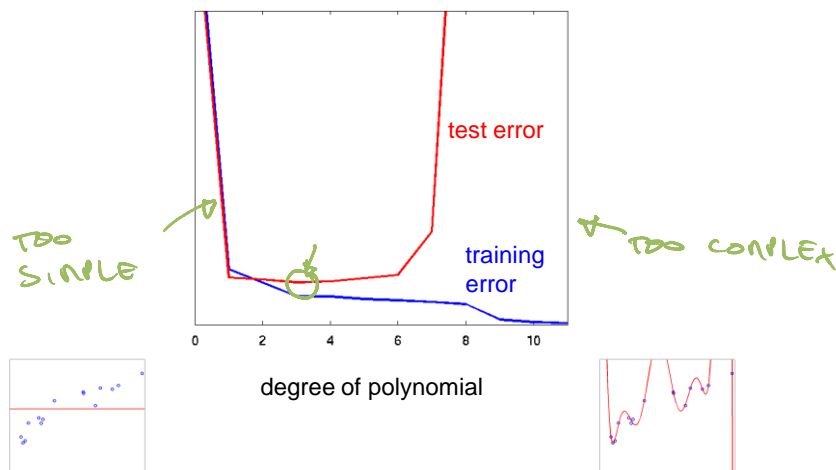
**Training error:**
overly optimistic

**Test error:**
approximation of **prediction error**
**as long as**
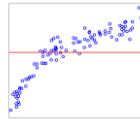**test data set never touched during training**

---

## Sweet spot for model complexity



TOO SINPLE

test error

training error

TOO CONPLEX

degree of polynomial

## Sweet spot for model complexity = Bias-variance tradeoff

**Bias:**
- faithfulness to the truth

**Variance:**
- sensitivity to randomness in training data



*HIGH BIAS*

*TRAINING SET OF SIZE 10*

*10th degree polynomial*

*LOW BIAS*

*HIGH VARIANCE*
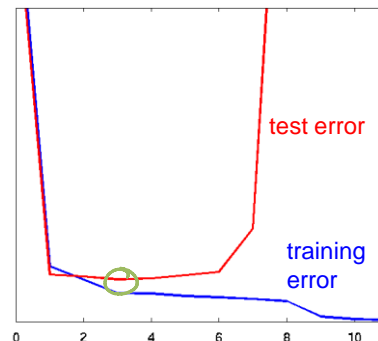
## true error

*DATA NOISE*

$$= bias^2 + variance + necessary\ evil$$

**Bias:**
- decreases with model complexity

**Variance:**
- increases with model complexity
- decreases with size of data set



test error

training error

degree of polynomial

*VARIANCE*

*BIAS*

**true error**

**= bias² + variance** + necessary evil

$X \sim p(x)$

$Y = f(X) + \varepsilon$ $\qquad \varepsilon \sim N(0, \sigma^2)$

prediction at $x_0$: $\qquad$ data $\rightarrow \hat{w}$ $\qquad$ RANDOM VAR

$\qquad$ RANDOM

$x_0 \mapsto \phi(x) \cdot \hat{w}$ $\qquad$ RANDOM VAR: $\hat{f}(x_0)$

error($x_0$) $= E_{Y,D}\left[ (Y - \hat{f}(x_0))^2 \mid X = x_0 \right]$

$E\left[ \left( (Y - E_Y[Y|x_0]) + (E_Y[Y|x_0] - E[\hat{f}(x_0)]) + (E_D[\hat{f}(x_0)] - \hat{f}) \right)^2 \mid x_0 \right]$

$f(x_0)$ $\qquad$ CONST.

$f(x_0)$

EXPECT = 0 $\qquad\qquad$ EXPECTATION = 0

---

**true error**

**= bias² + variance** + necessary evil

EXP=0 $\qquad\qquad$ CONST.

error($x_0$) $= E\left[ \left( (Y - f(X)) + (f(X) - E_D[\hat{f}(x_0)]) + (E_D[\hat{f}(x_0)] - \hat{f}(x_0)) \right)^2 \mid x_0 \right]$

$\varepsilon$ $\qquad\qquad$ EXP=0

$(A + B + C)^2 = A^2 + B^2 + C^2 + 2AB + 2AC + 2BC$

error($x_0$) $= E\left[ (Y - f(x_0))^2 \right] + (f(x_0) - E_D[\hat{f}(x_0)])^2$

$\sigma^2$ $\qquad\qquad$ BIAS²

error $= \int$ error$(x) p(x) dx$ $\qquad + E\left[ (\hat{f}(x_0) - E[\hat{f}(x_0)])^2 \right]$

VARIANCE

# Announcements

- project proposals due
  **this Wednesday at 10:30am**

- HW #5 OUT SOON
  (DUE OCT 7)

---

**Least squares fit**
**= max likelihood**
**for Gaussians**



Truth about salaries

$X \sim p(x)$   $N(0, \sigma^2)$

$Y \approx w \cdot \phi(x) + \varepsilon$

$p(Y | x) = N(w \cdot \phi(x), \sigma^2)$

$\max_w \prod_n p(Y_n | x_n)$  ← CONDITIONAL MLE

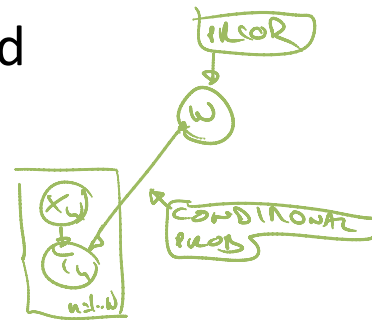$\max_w \sum_n \left( \text{const.} - \dfrac{(Y_n - w \cdot \phi(x_n))^2}{2\sigma^2} \right)$

$\min_w \sum_n (Y_n - w \cdot \phi(x_n))^2$

UNIFORM

# Beyond max likelihood for Gaussians

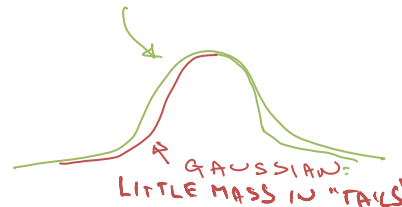- add a prior over **w**

  $\approx$ (MAP), prediction

- replace the Gaussian by a different model
  - different noise model
  - different support for **Y** $\in \mathbb{R}$

    $Y \in \{0, 1\}$

PRIOR

W

$x_N$

$y_N$

n=1..N

CONDITIONAL PROB

ALLOW OUTLIERS

FATTER TAILS

$p(\varepsilon) \propto \exp\{-|\varepsilon|\}$

GAUSSIAN: LITTLE MASS IN "TAILS"

---

# Regression with a Gaussian prior

$$Y \sim N(\phi(x) \cdot w, \, \omega^2)$$
$$w_j \sim N(0, \tau^2)$$

MAP:

$$\min_w \left[ -\log p(w | data) \right]$$

$$const + \frac{1}{2\omega^2} \sum_n (y_n - w \cdot \phi(x_n))^2 + \frac{1}{2\tau^2} \left( \sum_j w_j^2 \right) \quad / \, \omega^2$$

$$= \|w\|^2$$

NEG. LOG. LIKELIHOOD        NEG. LOG. PRIOR

$$\min_w \frac{1}{2} \|y - \Phi w\|^2 + \frac{1}{2} \left( \frac{\omega^2}{\tau^2} \right) \|w\|^2$$

ERROR        $\lambda$

CONTROL OVER COMPLEXITY

PUSHING TOWARDS $w = 0$

# Regression with a Gaussian prior

$Y \sim \phi(x) \cdot \widehat{\omega}$ + noise    *truth*

$w \sim N(0, \tau^2)$

$\hat{\omega} = (\phi^T \phi + \lambda I)^{-1} \phi^T y$

ALWAYS INVERTIBLE:

OPTIMIZATION
BETTER DEFINED
"REGULARIZED"
"
"SHRINKAGE"
"
NEG. LOG. PRIOR

$E[\hat{\omega}] = (\phi^T \phi + \lambda I)^{-1} \phi^T \phi \cdot w$

$A^{-1}$    $\lambda = 0$    $A$

$E[\hat{\omega}] = w$ if $\lambda = 0$: NO BIAS

$E[\hat{\omega}]$ is "shrinking" towards zero    INCREASING BIAS
as $\lambda$ increases

---

# Bias-variance tradeoff
# for ⟨ridge⟩ regression

REGULARIZATION
$\frac{\lambda}{2} \|w\|^2$

POLY OF
DEGREE 10



test error

training error

$10^5$    $10^{10}$

strength of regularization

BIAS
VARIANCE

# Bias-variance tradeoff
# for ridge regression