# Machine Learning Principles

## Class6 : Sept. 22

## Linear Regression III

## Instructor: Diana Kim

# Today's Lecture

1. Convex Optimization Theory

    - necessary & sufficient condition for optimality

    - equality constraint problem

    - inequality constraint problem

2. Regularization

    - three interpretations of MMSE with regularization

    - setting a right regularization parameter $\lambda*$: cross validation

3. Overfitting and Underfitting

- Optimization Theory:

Solving a convex optimization problem by using a Langrangian function

$$x^* = \min_{g(x) \le 0} f(x)$$

# [1] Local and Global Minimum (why optimization theory?)

In an ML problem, we need to solve an optimization problem, finding local / global minimum (suboptimal/optimal): MLE / MAP

- regression <u>with</u> constraint (regularization)

$$\arg\min_{\vec{w}} ||\vec{y} - \Phi \cdot \vec{w}||^2$$

- regression <u>without</u> constraint

$$\arg\min_{\vec{w}} ||\vec{y} - \Phi \cdot \vec{w}||^2$$
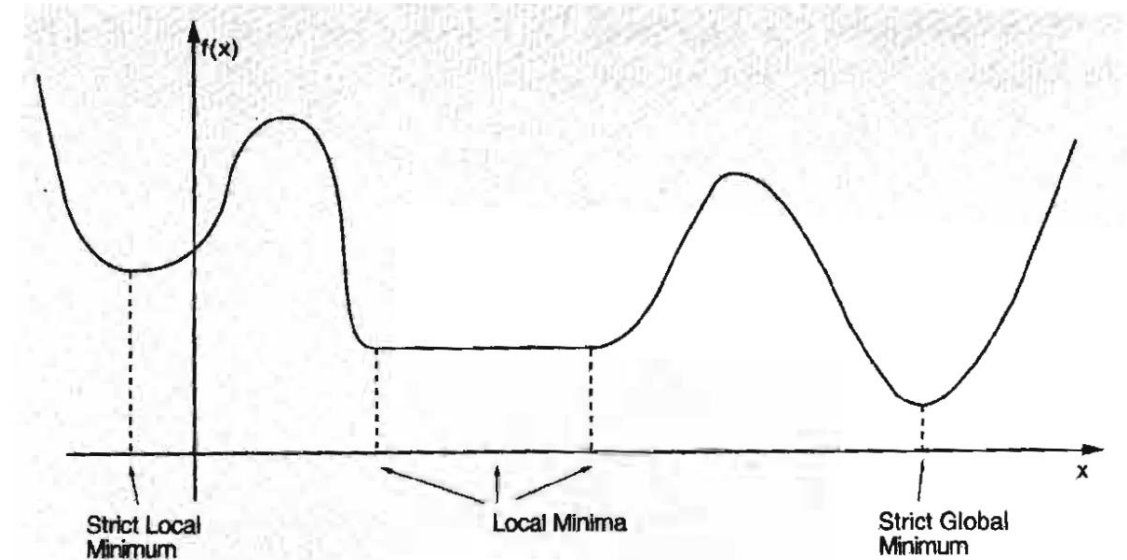$$\text{subject to} \quad ||\vec{w}||^2 \leq C$$

# [2] Local and Global Minimum

## Local Minimum $x^*$

$$f(x*) \leq f(x), \quad \exists \epsilon \quad s.t \quad ||x - x*|| < \epsilon \quad \forall x$$

## Global Minimum $x^*$

$$f(x*) \leq f(x) \quad \forall x$$

- By Taylor series

if $x^*$ is a local optimal,
then the Taylor approximation is non-negative:

$$f(x * + \Delta x) - f(x*) \approx \nabla f(x*)^t \Delta x + \frac{1}{2} \Delta x^t \nabla^2 f(x*) \Delta x \geq 0$$

# [5] Local and Global Minimum (Gradient Vector)

$$\nabla f(x) = \begin{bmatrix} \dfrac{\partial f(x)}{\partial x_1} \\[2ex] \dfrac{\partial f(x)}{\partial x_2} \\[2ex] \dfrac{\partial f(x)}{\partial x_3} \\[2ex] \dots \end{bmatrix}$$

**Gradient** of $f(x_1, x_2, x_3 \dots) = y \in \mathbb{R}$

- this is a vector.
- the collection of the first partial derivatives
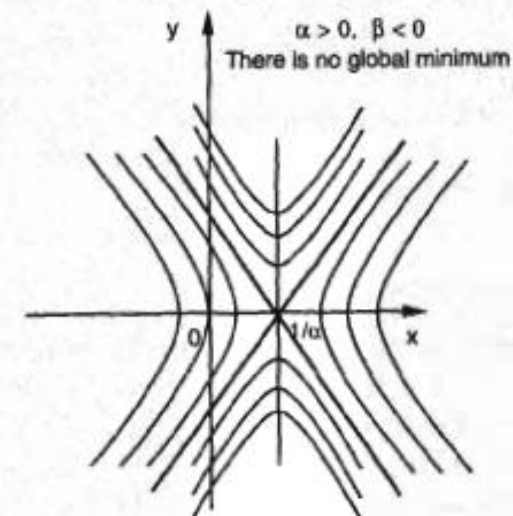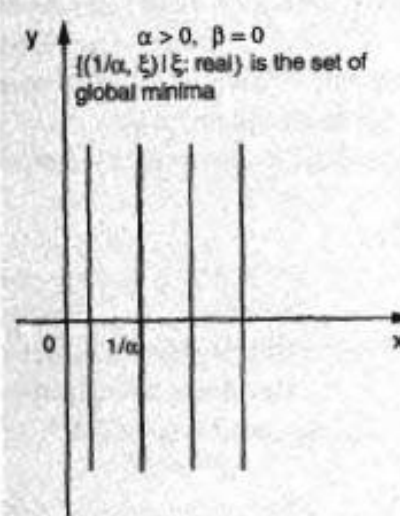- <span style="color:red">the direction of the greatest change of a scalar function</span>
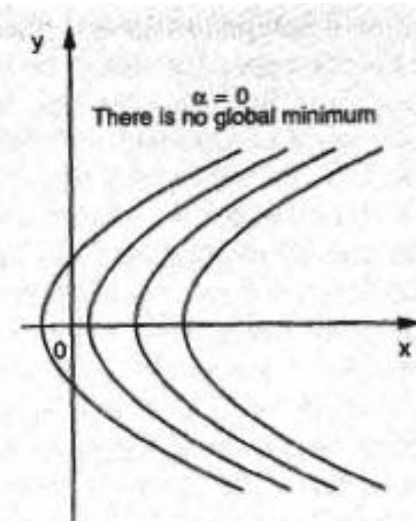
$$\nabla^2 f(x) = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial^2 x_1} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_3} & \cdots \\ \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \frac{\partial^2 f(x)}{\partial^2 x_2} & \frac{\partial^2 f(x)}{\partial x_2 \partial x_3} & \cdots \\ \frac{\partial^2 f(x)}{\partial x_3 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_3 \partial x_2} & \frac{\partial^2 f(x)}{\partial x_3 \partial x_3} & \cdots \\ \cdots & \cdots & \cdots & \cdots \end{bmatrix}$$

**Hessian Matrix** of $f(x_1, x_2, x_3 \dots) = y \in \mathbb{R}$
- this is a matrix.
- the collection of all second partial derivatives
- Concavity  of a function at the point $x^*$
- $\lambda_{min} \leq \Delta x^t \nabla^2 f(x *) \Delta x \leq \lambda_{max}$ where $||\Delta x||| = 1$
- The eigenvalues provide the concavity of principal axes (maximum/ minimum)

# [7] Local and Global Minimum (Hessian Matrix)

$$f(x, y) = \frac{1}{2} (\alpha x^2 + \beta y^2) - x$$

- **Equality Constraint Problem**

# [1] Equality Constraint Problem (example)

$$\min_{x} \quad f(x)$$

$$\text{s.t.} \quad h_i(x) = 0 \quad i = 1, ..., m$$

ex]
$$\min_{x} \quad x_1 + x_2$$

$$\text{s.t.} \quad x_1^2 + x_2^2 = 2$$

# [2] Equality Constraint Problem (necessary conditions)

$$\min_{x} \quad f(x)$$

$$\text{s.t.} \quad h_i(x) = 0 \quad i = 1, ..., m$$

- **Condition1**: let $x^*$ be a local minimum of $f$ s.t $h_i(x) = 0$ and $\nabla h_i(x^*)... \nabla h_i(x^*)$ are linearly independent. then, there exist a unique vector $\lambda^* = (\lambda_1^*, \lambda_2^*,.... \lambda_m^*)$ s.t

$$\nabla f(x*) + \sum_{i=1}^{m} \lambda_i^* \nabla h_i(x*) = 0$$

- **Condition2**: $y^t \{ \nabla^2 f(x*) + \sum_{i=1}^{n} \lambda_i^* \nabla^2 h_i(x*)) \} y \geq 0 \quad y \in V(x*)$

$$V(x*) = \{ y | \nabla h_i(x*)^t y = 0 \quad \forall i = 1, ..., m \}$$

# [3] Equality Constraint Problem (Lagrangian Multiplier Theorem)

$$\min_{x} \quad f(x)$$

$$\text{s.t.} \quad h_i(x) = 0 \quad i = 1, ..., m$$

- Condition1: let $x^*$ be a local minimum of f s.t $h_i(x) = 0$ and $\nabla h_i(x^*) ... \nabla h_i(x^*)$ are linearly independent. then, there exist a unique vector $\lambda^* = (\lambda_1^*, \lambda_2^*, ... \lambda_m^*)$ s.t

$$\nabla f(x*) + \sum_{i=1}^{m} \lambda_i^* \nabla h_i(x*) = 0$$

- Condition2: $\quad y^t \{ \nabla^2 f(x*) + \sum_{i=1}^{m} \lambda_i^* \nabla^2 h_i(x*)) \} y \geq 0 \quad y \in V(x*)$

$$V(x*) = \{ y | \nabla h_i(x*)^t y = 0 \quad \forall i = 1, ..., m \}$$

Q: What if we define a new function? $L(x, \lambda) = f(x) + \sum_{i=1}^{m} \lambda_i h_i(x)$

CS 461: class #6

# [4] Equality Constraint Problem (Lagrangian function)

- Lagrangian function/ unconstrained function

$$L(x, \lambda) = f(x) + \sum_{i=1}^{m} \lambda_i \, h_i(x)$$

- two necessary optimality conditions for $L(x, \lambda)$

$$\nabla f(x*) + \sum_{i=1}^{m} \lambda_i^* \nabla h_i(x*) = 0 \qquad h_i(x*) = 0 \quad \forall i = 1, 2..., m$$

$$y^t \{ \nabla^2 f(x*) + \sum_{i=1}^{m} \lambda_i^* \nabla^2 h_i(x*)) \} y \geq 0 \quad y \in V(x*)$$

$$V(x*) = \{ y | \nabla h_i(x*)^t y = 0 \quad \forall i = 1, ..., m \}$$

# [5] Equality Constraint Problem (Lagrangian function)

[Lagrangian Function/ unconstrained function]

$$L(x, \lambda) = f(x) + \sum_{i=1}^{m} \lambda_i \, h_i(x)$$

The necessary conditions for optimal points for the unconstrained function (Lagrangain) is equivalent to the constrained primary problem.
Therefore, we could fine optimal solutions for the primal problem by solving the necessary conditions on Lagrangian function $L(x, \lambda)$

# [6] Equality Constraint Problem (Lagrangian example)

Consider the problem

$$\text{minimize} \quad \tfrac{1}{2}\left(x_1^2 + x_2^2 + x_3^2\right)$$

$$\text{subject to} \quad x_1 + x_2 + x_3 = -3.$$

Q: Lagrangian function?

- Inequality Constraint Problem

# [1] Inequality Constraint Problem (necessary conditions)

$$\min_{x} \quad f(x)$$

$$\text{s.t.} \quad g_i(x) \leq 0 \quad i = 1, ..., m$$

\* pay attention! the direction of inequality!

# [2] Inequality Constraint Problem (necessary conditions)

$$\min_{x} \quad f(x)$$

$$\text{s.t.} \quad g_i(x) \leq 0 \quad i = 1, ..., m$$

- two possible cases for $x^*$

  (1) $x^*$ inside of $g_i(x) < 0$

  (2) $x^*$ on the boundary of $g_i(x) = 0$

$X$

# [3] Inequality Constraint Problem (necessary conditions)

$$\min_{x} \quad f(x)$$

$$\text{s.t.} \quad g_i(x) \le 0 \quad i = 1, ..., m$$

▪ two possible cases and their necessary conditions

(1) $x^*$ inside of the manifold by $g_j(x) < 0$

$$\rightarrow \nabla f(x*) = 0$$

(2) $x^*$ on the boundary of $g_h(x) = 0$ there exist $\lambda_h$

$X$

$$\rightarrow \nabla f(x*) + \lambda_h \nabla g_h(x*) = 0 \quad \text{Q: sign } \lambda_h \text{ ?}$$

# [4] Inequality Constraint Problem (necessary conditions)

# [5] Inequality Constraint Problem (KKT necessary conditions)

Let $x^*$ be a local minimum of the problem          [**K**arush–**K**uhn–**T**ucker conditions]

$$\min_x \quad f(x)$$

$$\text{s.t.} \quad g_i(x) \leq 0 \quad i = 1, ..., m$$

important to set the inequality this form (less than or equal to)!

Then, there exist $\lambda_i$, $i = 1, ..., m$ such that

- stationary condition

(1)    $\nabla f(x*) + \sum_{i=1}^{m} \lambda_i \nabla g_m(x*) = 0$

- complementary slackness condition
  $\lambda_i \cdot g_i(x^*) = 0$

(2)    $\begin{cases} \lambda_j \geq 0 & j = 1, ..., r \\ \lambda_j = 0 & \forall j \notin A(x*) \end{cases}$    $A(x^*)$ is the set of active constraints at $x^*$

# [6] Inequality Constraint Problem (KKT necessary conditions)

Let $x^*$ be a local minimum of the problem

$$\min_{x} \quad f(x)$$

$$\text{s.t.} \quad g_i(x) \leq 0 \quad i = 1, \dots, m$$

Then, there exist $\lambda_i, \; i = 1, \dots, m$ such that

(1) $\nabla f(x*) + \sum_{i=1}^{m} \lambda_i \nabla g_m(x*) = 0$

(2) $\begin{cases} \lambda_j \geq 0 & j = 1, \dots, r \\ \lambda_j = 0 & \forall j \notin A(x*) \end{cases}$

(3) $g(x^*) \leq 0$    ▪ <span style="color:red">Primary feasibility</span>

# [7] Inequality Constraint Problem (example 1)

Consider the problem

$$\text{minimize} \quad \tfrac{1}{2}\left(x_1^2 + x_2^2 + x_3^2\right)$$

$$\text{subject to} \quad x_1 + x_2 + x_3 \leq -3.$$

Then for a local minimum $x^*$, the first order necessary condition [cf. Eq. (3.47)] yields

$$x_1^* + \mu^* = 0,$$

$$x_2^* + \mu^* = 0,$$

$$x_3^* + \mu^* = 0.$$

From Nonlinear Programming, Bertsekas Example 3.3.1

- Regularization as an optimization problem

$$\arg\min_{\vec{w}} ||\vec{y} - \Phi \cdot \vec{w}||^2 + \lambda^*(||\vec{w}||^2)$$

# [1] Regularization by an optimization problem

In a ML problem, we need to solve an optimization problem, finding local / global minimum (suboptimal/optimal).

- regression <u>without</u> constraint

$$\arg\min_{\vec{w}} ||\vec{y} - \Phi \cdot \vec{w}||^2$$

- regression <u>with</u> constraint (regularization)

$$\arg\min_{\vec{w}} ||\vec{y} - \Phi \cdot \vec{w}||^2$$
$$\text{subject to} \quad ||\vec{w}||^2 \leq C$$

# [2] Regularization by an optimization problem (Lagrangian form)

$$\arg\min_{\vec{w}} ||\vec{y} - \Phi \cdot \vec{w}||^2$$
$$\text{subject to} \quad ||\vec{w}||^2 \leq C$$

$\longleftrightarrow$

$$\arg\min_{\vec{w}} ||\vec{y} - \Phi \cdot \vec{w}||^2 + \lambda^*(||\vec{w}||^2 - C)$$

- $C \propto \dfrac{1}{\lambda*}$ "but does not need to know the exact value of C"

- an example problem will be discussed in recitation #3 Sept. 24[th].

# [3] Regularization by an optimization problem (Lagrangian form)

$$\underset{\vec{w}}{\arg\min} \, ||\vec{y} - \Phi \cdot \vec{w}||^2$$
$$\text{subject to} \quad ||\vec{w}||^2 \leq C$$

$$\longleftrightarrow$$

$$\underset{\vec{w}}{\arg\min} \, ||\vec{y} - \Phi \cdot \vec{w}||^2 + \lambda^*(||\vec{w}||^2 - C)$$

- according to $C$ we define,
  optimal Lagrangian $\lambda *$ will be different!

- constant addition/subtraction won't change $x^*$

$$\underset{\vec{w}}{\arg\min} \, ||\vec{y} - \Phi \cdot \vec{w}||^2 - (\lambda^* C) + \lambda^*(||\vec{w}||^2)$$

# [4] Regularization by an optimization problem (Lagrangian form)

$$\arg\min_{\vec{w}} ||\vec{y} - \Phi \cdot \vec{w}||^2$$

$$\text{subject to} \quad ||\vec{w}||^2 \leq C$$

$$\arg\min_{\vec{w}} ||\vec{y} - \Phi \cdot \vec{w}||^2 + \lambda^*(||\vec{w}||^2)$$

- in regularized regression learning, we will change $\lambda^*$ and test its performance to find a good $\lambda^*$ (empirically)

# [5] Regularization by an optimization problem (Ridge & Lasso)

$$\arg\min_{\vec{w}} ||\vec{y} - \Phi \cdot \vec{w}||^2 + \lambda^*(||\vec{w}||^2) \quad \text{[Ridge regularization]}$$

$$\arg\min_{\vec{w}} ||\vec{y} - \Phi \cdot \vec{w}||^2 + \lambda^*(||\vec{w}||) \quad \text{[Lasso regularization]}$$

Ridge    Lasso

- the constraints regulate the magnitude of $w$ (parameters), the model complexity. Lasso gives a sparse solution.

From Bishop Chap Figure 3.4

- How could we set $\lambda *$? (controlling C)

# [1] selection of $\lambda *$

- we have data points $\{(x_1, y_1), (x_2, y_2), (x_3, y_3), ..., (x_N, y_N)\}$
- split the data into train/ test
- set basis functions (polynomial / Gaussian)

# [2] selection of $\lambda*$

- we have data points $\{(x_1, y_1), (x_2, y_2), (x_3, y_3), ..., (x_N, y_N)\}$
- split the data into train/ validation/ test
- set hypothetical basis functions (polynomial / Gaussian)

[train]

- PCA of train data for dim-reduction and whitening.
  (save the PCA blocks computed using training set)
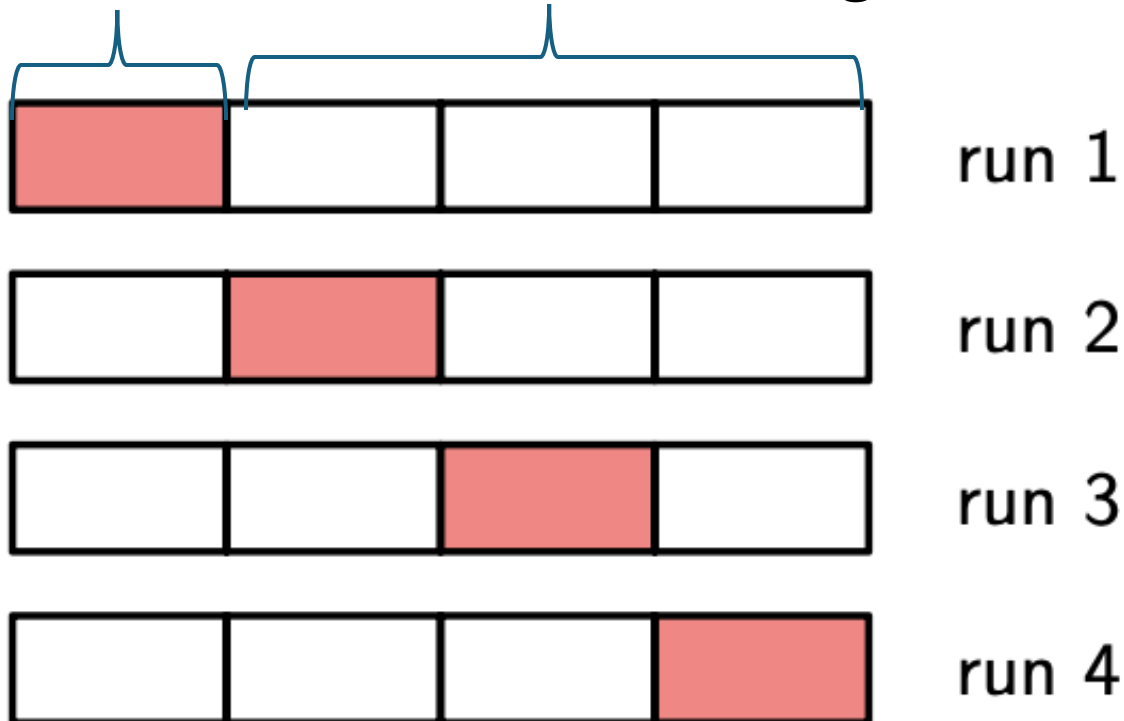- set the normal equations for the various $\lambda*$ to estimate $w$

[validation]

- PCA of validation data and whitening with the same block in training.
- compute error between the ground truth $y$ and prediction $y'$: $||y - y'||^2$
- choose the best $\lambda*$

# [3] selection of $\lambda *$ (validation set)

Validation set is a separate hold out set.
One example ratio = train: validation: test = 80: 10 : 10

- <u>to select regularization parameters</u>
- to finalize basis set
- to finalize feature map (# features dim)

# [**recall] Linear Regression by MAP (solving the optimization problem)

$$\nabla J(W) = -2\Phi^t \cdot \vec{Y} + 2\Phi^t \cdot \Phi \cdot \vec{W} + 2\lambda * \cdot \vec{W} = 0$$

$$\leftrightarrow \Phi^t \cdot \Phi \cdot \vec{W} + \lambda * \cdot \vec{W} = \Phi^t \cdot \vec{Y}$$

$$\leftrightarrow V \begin{bmatrix} \lambda_1 + \lambda* & 0 & ... & 0 \\ 0 & \lambda_2 + \lambda* & ... & 0 \\ ... & ... & ... & ... \\ 0 & ... & ... & \lambda_m + \lambda* \end{bmatrix} V^t \cdot \vec{W} = V^t \lambda^{1/2} E^t \vec{Y}$$

$$\leftrightarrow \vec{W} = V \begin{bmatrix} \dfrac{1}{\lambda_1 + \lambda*} & 0 & ... & 0 \\ 0 & \dfrac{1}{\lambda_2 + \lambda*} & ... & 0 \\ ... & ... & ... & ... \\ 0 & ... & ... & \dfrac{1}{\lambda_m + \lambda*} \end{bmatrix} E^t \vec{Y}$$

[by the $\lambda^*$ we can avoid the case parameters gets too large.]

# [4] selection of $\lambda *$

[train]

- PCA of train data for dim-reduction and whitening. (save the PCA blocks computed using training set)
- set normal equation with various $\lambda *$ to estimate $w$

[validation]

- PCA of validation data and whitening with the same block in training.
- compute error between the groud turth $y$ and prediction $y'$: $||y - y'||^2$
- Choose the best $\lambda *$

[test: report the final test results!]

- PCA of test data and whitening with the same block in
- compute error between the groud turth $y$ and prediction $y'$: $||y - y'||^2$

# [5] selection of $\lambda *$ (cross validation)

hold out  used in training



run 1

run 2

run 3

run 4

$$L = \frac{1}{S} \sum_{i=1}^{S} L_s$$

The example for the case $S = 4$

- **Overfitting & Underfitting**
(variance vs. bias)

effective feature engineering is the key to linear regression problem.

- choice of basis functions: nature of target tasks
- # (num) features: complexity but need to consider # data points
- no collinearity.

The higher polynomials are expressive, but we should be careful when using them when with a limited number of data points.

# [1] Overfitting & Underfitting

<span style="color:red">Given a set of basis functions and data,</span>
we can have the two problematic cases (if having inappropriate complexity):

(1) **Overfitting** / (2) **Underfitting**

# [2] Overfitting & Underfitting (example of overfitting)



original

original

# data = 10
noise $\sigma \approx \sqrt{0.15}$

(1) regression with P0-p9 with #10 data points
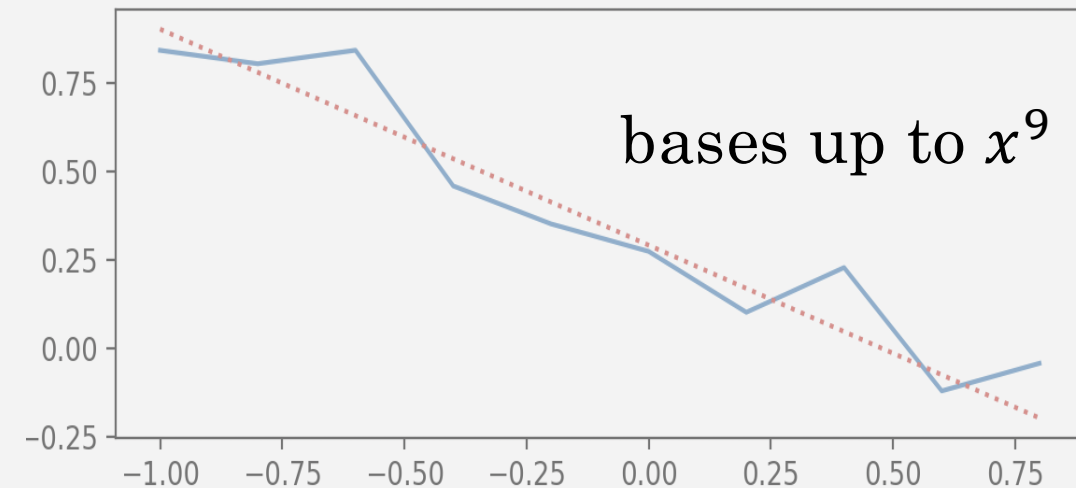
(1) regression with P0-p2 with #10 data points

bases up to $x^9$

bases up to $x^2$

original

original

(1) regression with P0-p9 with #10 data points

(1) regression with P0-p2 with #10 data points

bases up to $x^9$

bases up to $x^2$

+ the left-side high complexity model captured the noise in data.

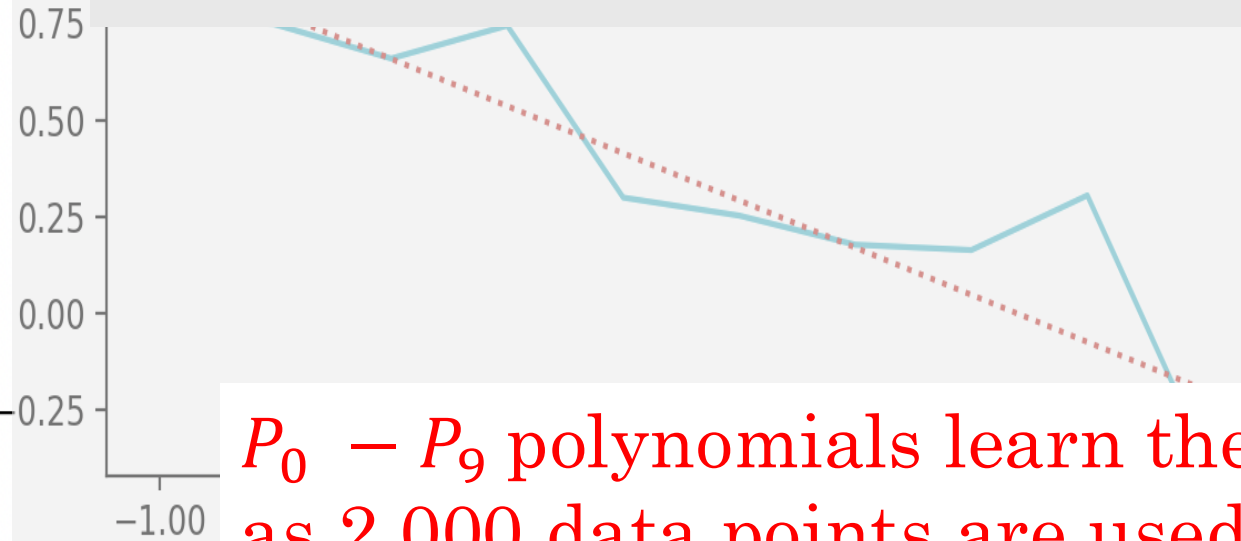# [4] Overfitting & Underfitting (overfitting: too complex)

- overfitting

when model complexity is too high relatively to # training data,
then the complex model fits to noise/ undefined by data. one phenomenon is
the large gap between train vs. test; we say the model has a high variance.
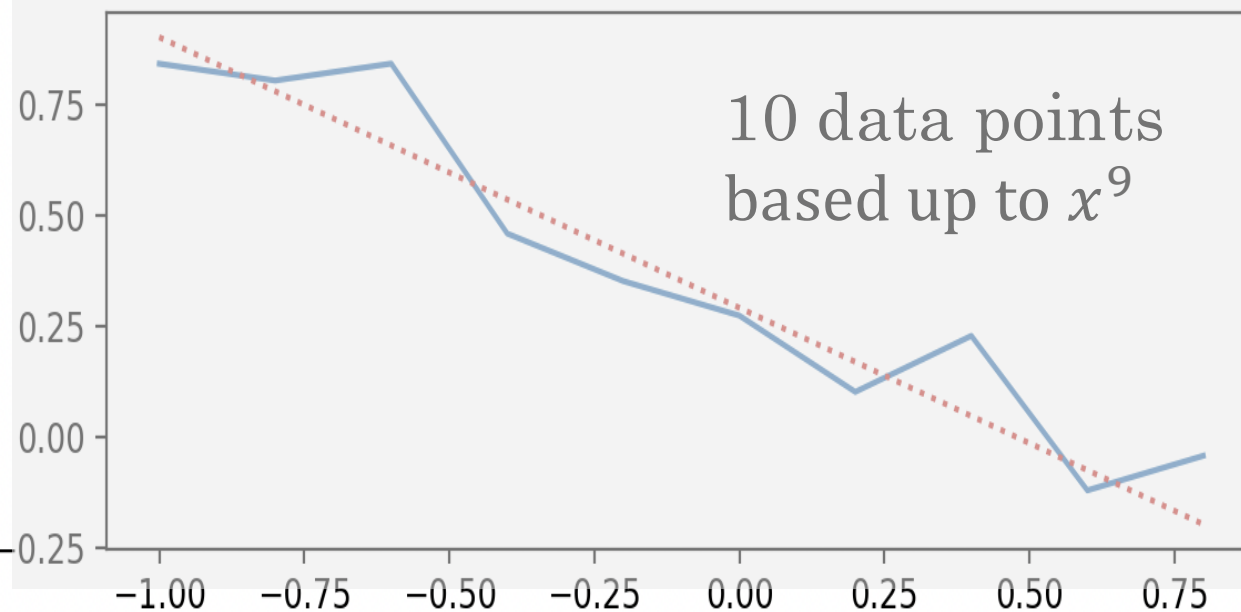
Q: what if we have enough training data?

# [5] Overfitting & Underfitting (overfitting: too complex)

- overfitting

when model complexity is too high relatively to # training data,
then the complex model fits to noise/ undefined by data. one phenomenon is
the large gap between train vs. test; we say the model has a high variance.

Q: what if we have enough training data?
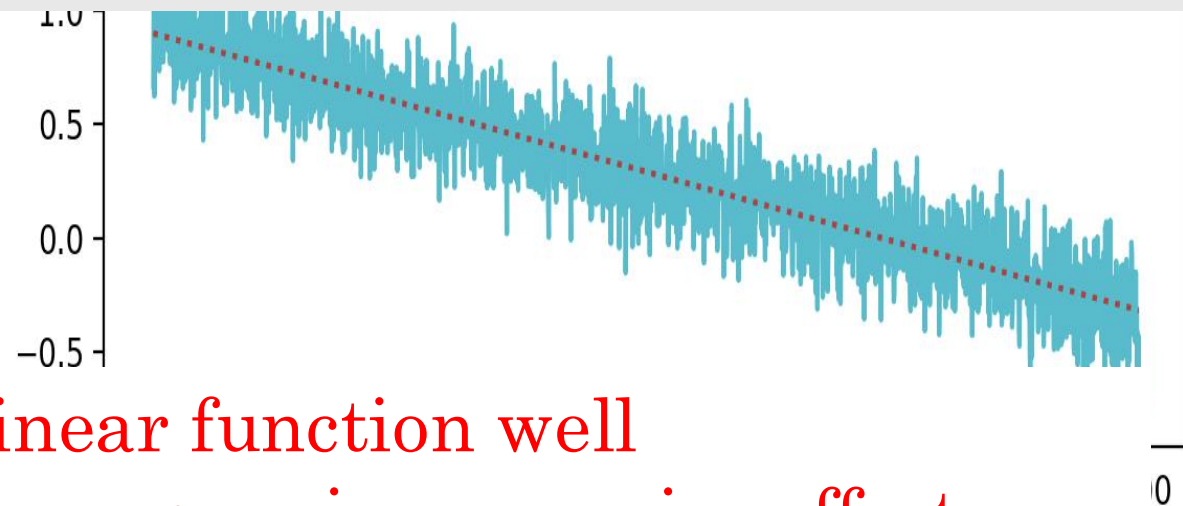    $\sum_{i=1}^{\infty} \varepsilon_i = 0$  noise effect diminishes.

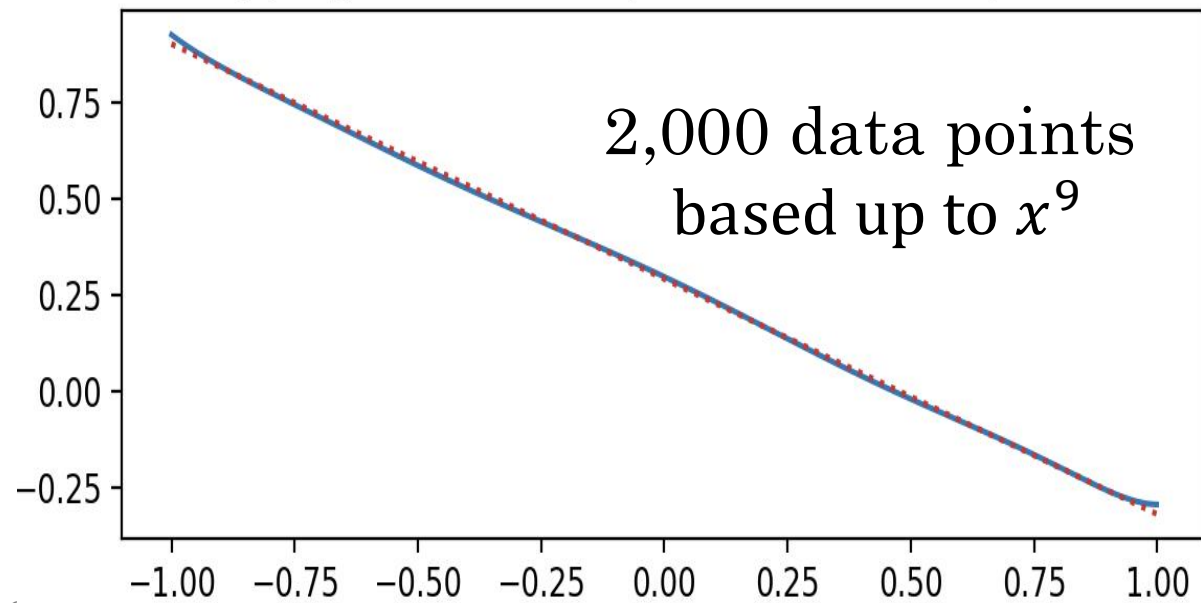$P_0 - P_9$ polynomials learn the linear function well as 2,000 data points are used for regression. no noise effect.

(1) regression with P0-p9 with #10 data points

(1) regression with P0-p9 with #2000 data points

10 data points based up to $x^9$

2,000 data points based up to $x^9$

# [7] Overfitting & Underfitting (with large amount of data)

[from Bishop p 9.]

<span style="color:red">the larger data data set,
the more complex the model that we can afford to fit to the data.</span>
One rough heuristic that is sometimes advocated that the number of data points should be no less than some multiple of the number of adaptive parameters.
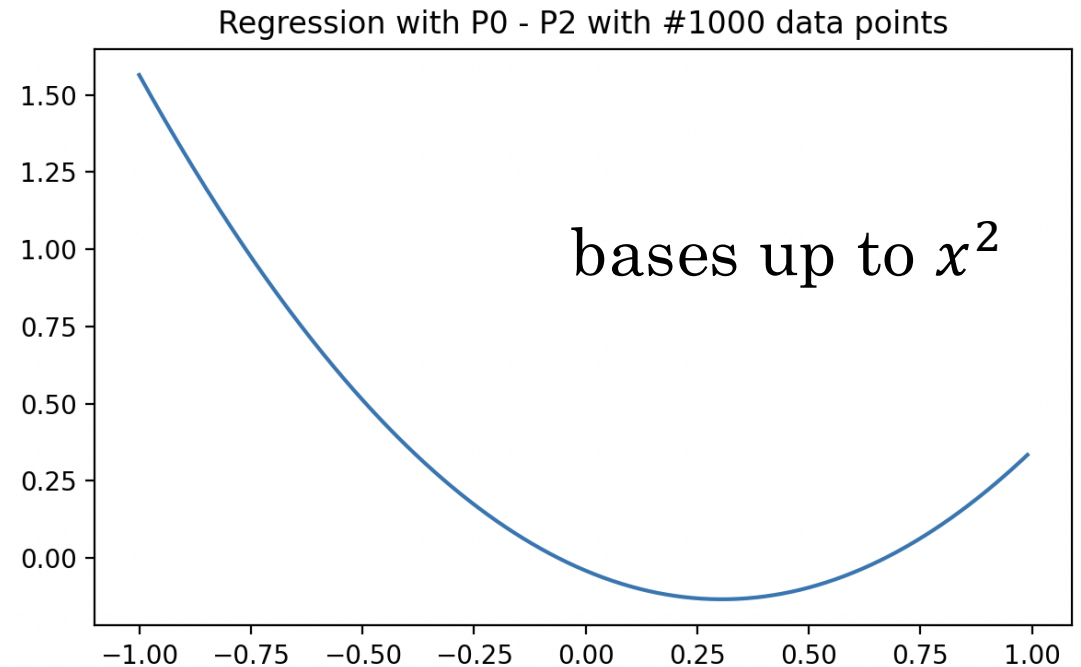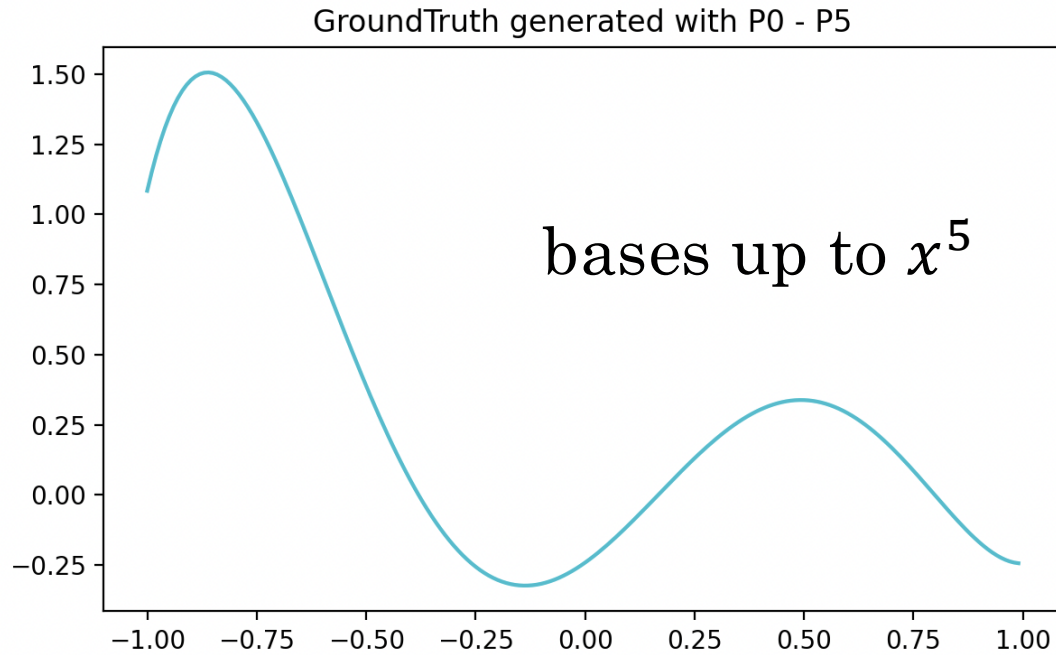
# [8] Overfitting & Underfitting

Q: how could we prevent the overfitting as we have limited data?
   do we need to collect more data points?
   do we need to reduce the number of features?
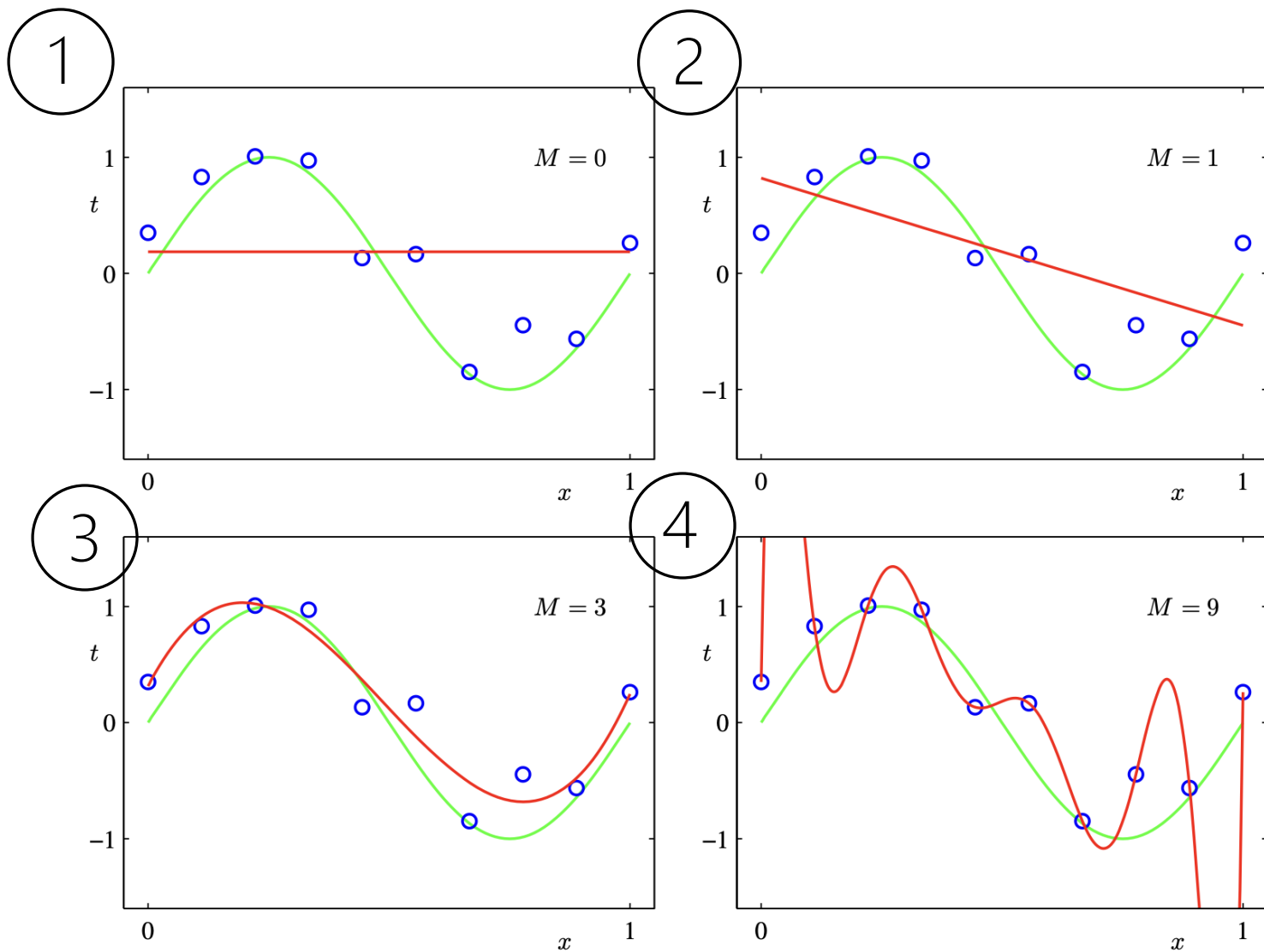
# [9] Overfitting & Underfitting (underfitting)

when model complexity
is not enough for the ground truth model,
no way to learn even when we have enough data. (no hope to learn)

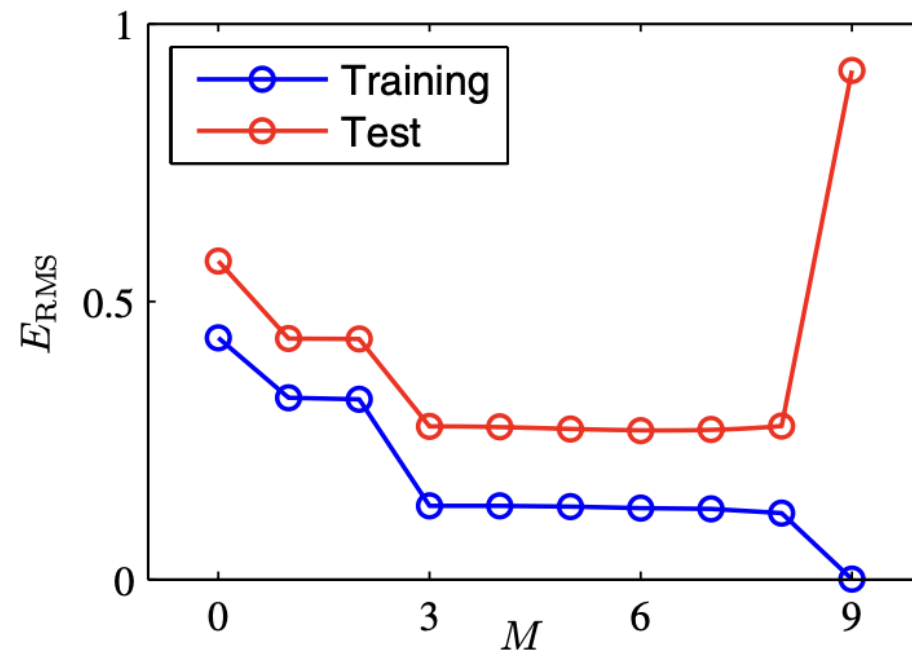# [10] Overfitting & Underfitting (example of underfitting)



GroundTruth generated with P0 - P5

bases up to $x^5$

Regression with P0 - P2 with #1000 data points

bases up to $x^2$

- underfitting example:
  # data is 1,000 but  as the basis functions are limited by $P_0 - P_2$
  no way to capture the original function generated by $P_0 - P_5$

# [11] Overfitting & Underfitting (textbook examples)



[from Bishop Figure 1.4 and 1.5]

Q: which one shows overfitting/ underfitting?

## [12] Overfitting & Underfitting (symptoms)

- **overfitting**
  train performance good but test performance is bad
  the performance gap between test and train set (<span style="color:red">poor generalization/ high variance</span>)

- **underfitting**
  both train and test performances are bad (<span style="color:red">biased/ high bias</span>)

Q: how could we learn a model in right complexity?
    for the primary goal in ML, high generalization!

- The Bias Variance Decomposition

test error $=$ variance $+ bias^2 +$ intrinsic error

# [1] Bias & Variance Decomposition (Empirical MSE & Expected MSE)

- suppose we had a data set: $D = \{(\vec{x_i}, t_i),\ i=1,2,...n\}$
- we had a set of basis function
- we learned a regression model $y(\vec{x})$ from $D$.

# [2] Bias & Variance Decomposition (Empirical MSE & Expected MSE)

- suppose we had a data set: $D = \{(\vec{x_i}, t_i),\ i=1,2,...N\}$
- we had a set of basis functions
- we learned a regression model $y(\vec{x}; D)$ from $D$.

(1) empirical MSE error (this is the one we usually do.)

$$L = \frac{1}{N} \sum_i^N \{y(x_i; D) - t_i\}^2$$

(2) expected MSE error, assuming we know the density $f(\vec{x}, t) = f(t|\vec{x})\, f(\vec{x})$

$$E[L] = \int_x \int_t \{y(x; D) - t)\}^2 f(t|x) f(x)\, \mathrm{d}t\, \mathrm{d}x$$

# [3] Bias & Variance Decomposition (decomposition)

$$E[L] = \int_x \int_t \{y(x; D) - t)\}^2 f(t|x) f(x) \, dt \, dx$$

$$E[L] = \int_x \int_t (y(x; D) - h(x) + h(x) - t)^2 f(t|x) f(x) \, dt \, dx$$

$$= \int_x \int_t (y(x; D) - h(x))^2 + (h(x) - t)^2 + 2(y(x; D) - h(x))(h(x) - t) \, dt \, dx$$

- $h(x)$ is an optimal function minimizing the modeling error given density $f(x, t)$

$$h(x) = \arg\min_{y(x)} \int_x \int_t \{y(x) - t)\}^2 f(t|x) f(x) \, dt \, dx$$

# [4] Bias & Variance Decomposition (optimal $h(x) = E[T|x]$)

Suppose y(x) is our model, data follows the density $f(t, x) = f(t|x) \cdot f(x)$

$$E[L] = \int_x \int_t (y(x) - t)^2 f(t|x) f(x) \, dt \, dx$$

$$J(y^*(x)) = \int_t (y^*(x) - t)^2 f(t|x) \, dt$$

$$\frac{\partial J}{\partial y^*(x)} = \int_t (2y^*(x) - 2 \cdot t) f(t|x) \, dt = 0$$

$$y^*(x) = \int_t t f(t|x) \, dt = \boxed{E[T|x]}$$

**for $f(x) \geq 0$, we can set up minimization only for $\int_t (y(x) - t)^2 f(t|x)$

$E[T|x]$ is a function of $x$

# [5] Bias & Variance Decomposition (decomposition)

$$E[L] = \int_x \int_t \{y(x; D) - t)\}^2 f(t|x) f(x) \, \mathrm{d}t \, \mathrm{d}x$$

$$E[L] = \int_x \int_t (y(x; D) - h(x) + h(x) - t)^2 f(t|x) f(x) \, \mathrm{d}t \, \mathrm{d}x$$

$$= \int_x \int_t (y(x; D) - h(x))^2 + (h(x) - t)^2 + 2(y(x; D) - h(x))(h(x) - t) \, \mathrm{d}t \, \mathrm{d}x$$

? $\qquad$ $\epsilon$ $\qquad$ zero

$$t = h(x) + \epsilon$$

$$\int_x \int_t 2(y(x; D) - h(x))(h(x) - t) f(t|x) f(x) \, dt \, dx$$

$$\int_x 2(y(x; D) - h(x)) \int_t (h(x) - t) \, \mathrm{d}t f(t|x) \, \mathrm{d}t f(x) \, \mathrm{d}x$$

$$\int_x 2(y(x; D) - h(x))(\underline{E[T|x] - E[T|x]}) f(x) \, \mathrm{d}x = 0$$

zero

# [6] Bias & Variance Decomposition (decomposition)

$$E[L] = \int_x \int_t \{y(x; D) - t)\}^2 f(t|x) f(x) \, dt \, dx$$

$$E[L] = \int_x \int_t (y(x; D) - h(x) + h(x) - t)^2 f(t|x) f(x) \, dt \, dx$$

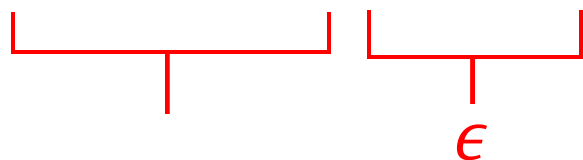$$= \int_x \int_t (y(x; D) - h(x))^2 + (h(x) - t)^2 + 2(y(x; D) - h(x))(h(x) - t) \, dt \, dx$$

$\epsilon$

$$\int_x \int_t (y(x; D) - h(x))^2 f(t|x) f(x) \, dt \, dx$$

$$\int_x (y(x; D) - E_D[y(x; D)] + E_D[y(x; D)] - h(x))^2 f(x) \, dx$$

expectation over different data sets

$$\int_x (y(x; D) - E_D[y(x; D)])^2 + (E_D[y(x; D)] - h(x))^2 f(x) \, dx$$

variance          bias

# [7] Bias & Variance Decomposition (decomposition)

- suppose we had a data set: $D = \{(\vec{x_i}, t_i), i=1, 2, \ldots, N\}$
- we had a set of basis functions
- we learned a regression model $y(\vec{x}; D)$ from $D$.
- the expected error is decomposed into three terms
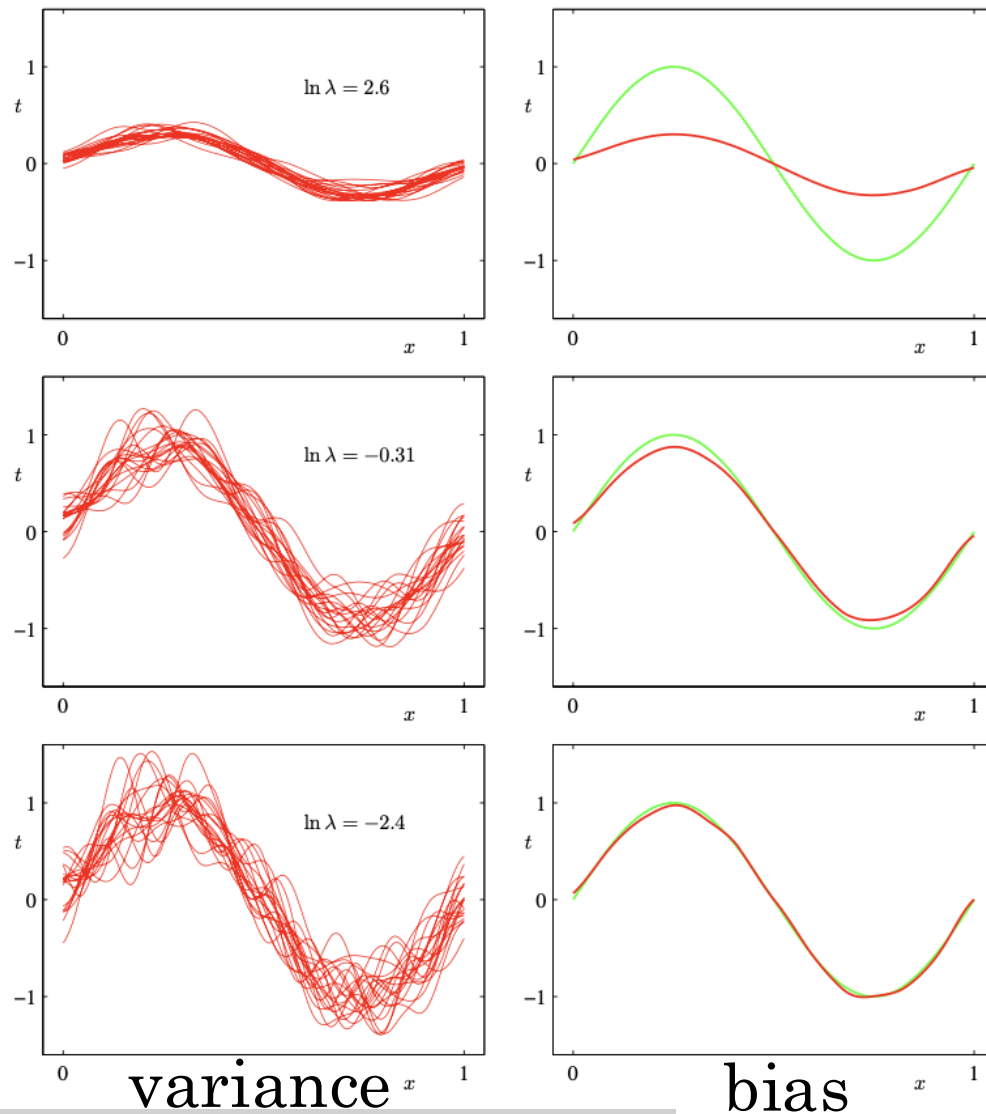  $E[L]$ = variance + $bias^2$ + intrinsic error

- Intrinsic Error: $\int_x \int_t (E[T|x] - t)^2 f(t|x) f(x) \, dt \, dx$

- Variance: $\int_x (y(x; D) - E_D[y(x; D)])^2 \, dx$

- Bias: $\int_x \{E_D[y(x; D)] - E[T|x]\}^2 f(x) \, dx$

variance and bias vary depending on model complexity!
(basis set or regularization parameter)

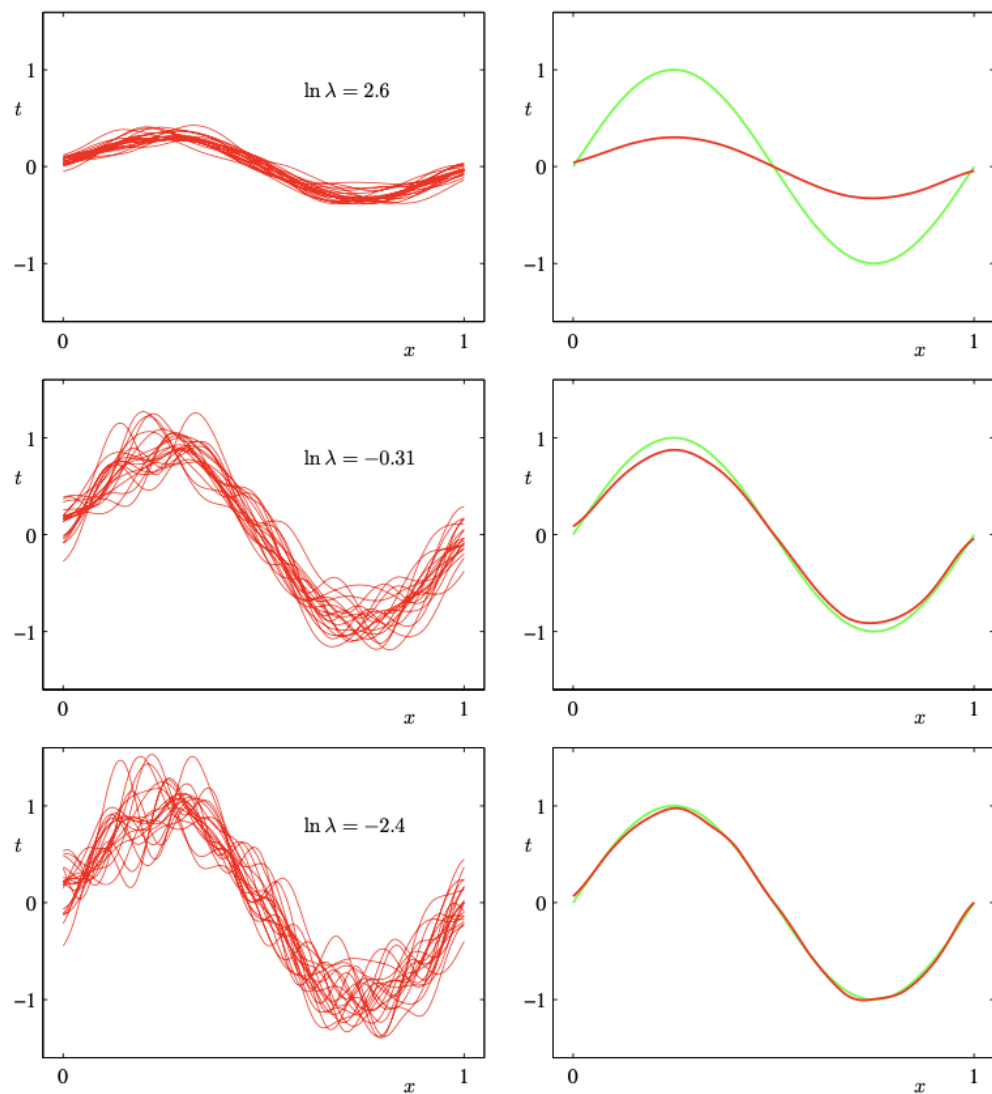# [8] Bias & Variance Decomposition (trade off textbook example)
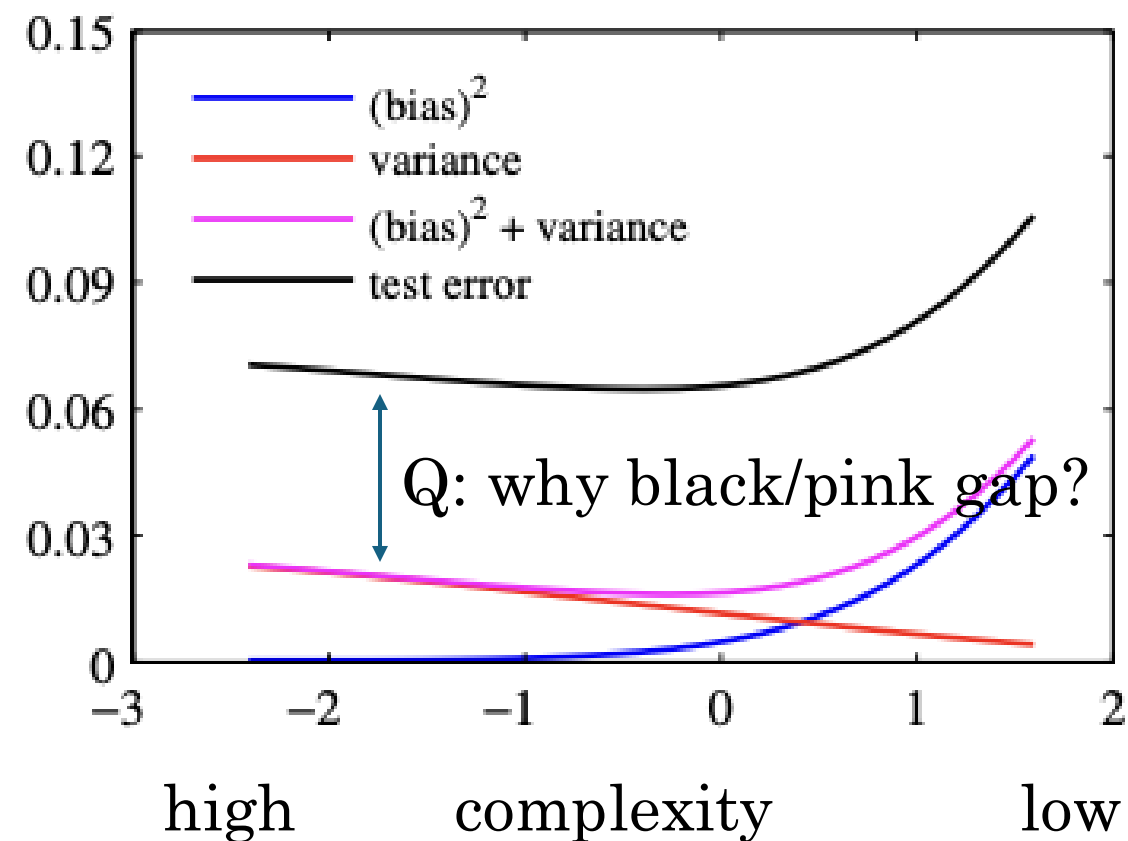


[from Bishop p. 150]

increasing complexity

variance    bias

Q: how changes the variance & bias?

# [9] Bias & Variance Decomposition (trade off textbook example)



[from Bishop p. 150]

increasing complexity

$\ln \lambda = 2.6$

$\ln \lambda = -0.31$

$\ln \lambda = -2.4$

- (bias)$^2$
- variance
- (bias)$^2$ + variance
- test error

Q: why black/pink gap?

high          complexity          low

- Controlling Effective Model Complexity

  (strategy to avoid over/underfitting)

## [1] Controlling Effective Model Complexity (strategies)

- **underfitting**
  both  train and test performances are bad (biased/ high bias)

- **overfitting**
  train performance good but test performance is bad
  the performance gap between test and train set (poor generalization/ high variance)

Q: how could we learn a model in right complexity?
     for the primary goal in ML, high generalization!

# [2] Controlling Effective Model Complexity (colliding strategies)

- **underfitting**
  both train and test performances are bad (biased/ high bias)
  increasing the order of model hypothesis space

- **overfitting**
  train performance good but test performance is bad
  the performance gap between test and train set (poor generalization/ high variance)
  decreasing the order of model hypothesis space

  Q: what will be the best strategies?

# [2] Controlling Effective Model Complexity (colliding strategies)

- **Underfitting**
  both train and test performances are bad (biased/ high bias)
  increasing the order of model hypothesis space

- **Overfitting**
  train performance good but test performance is bad
  the performance gap between test and train set (poor generalization/ high variance)
  decreasing the order of model hypothesis space

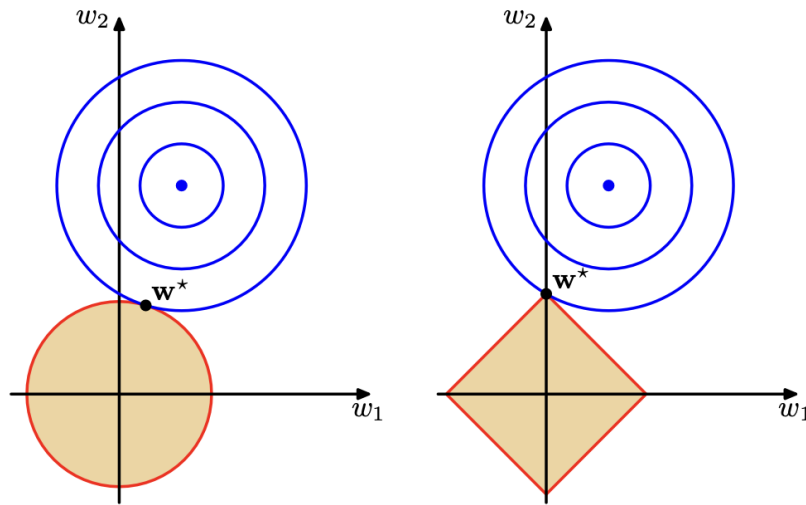Q: what will be the best strategies?
  collect many data samples but expansive. (ideal solution but not practical)
  then, how do we control model complexity?

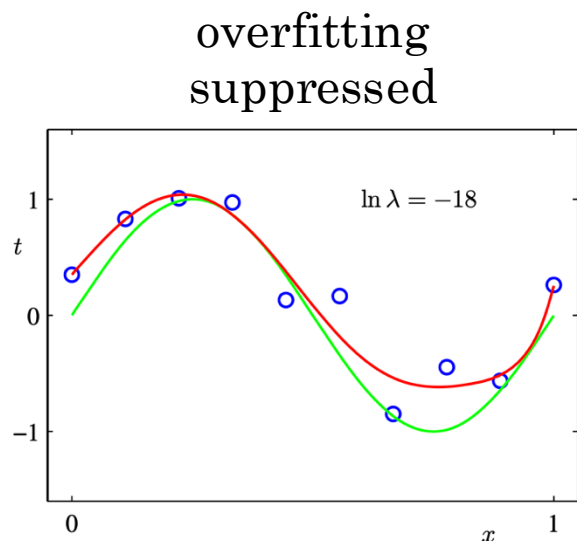# [3] Controlling Effective Model Complexity (regularization)

$$\arg\min_{\vec{w}} ||\vec{y} - \Phi \cdot \vec{w}||^2 + \lambda^*(||\vec{w}||^2)$$   [Ridge regularization]

$$\arg\min_{\vec{w}} ||\vec{y} - \Phi \cdot \vec{w}||^2 + \lambda^*(||\vec{w}||)$$   [Lasso regularization]
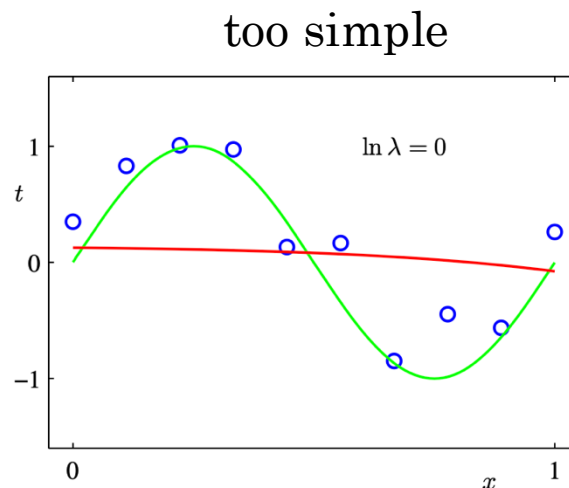


- In Bayesian model, the effective number of parameters adapts automatically to the size of the data set.

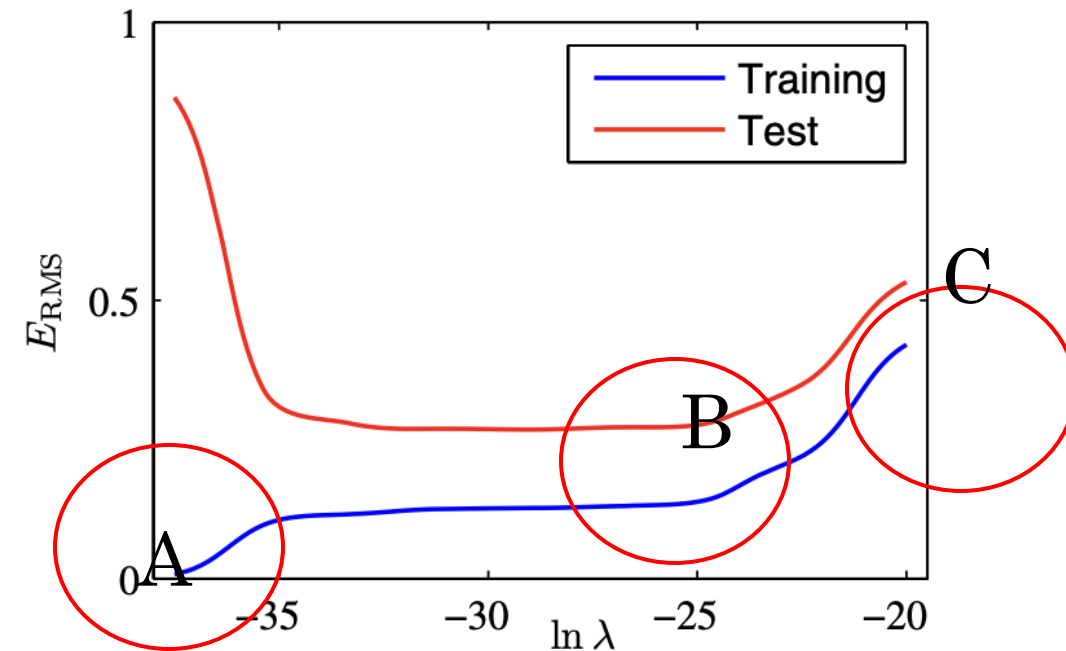# [4] Controlling Effective Model Complexity (regularization by lambda)

overfitting
suppressed

$\ln \lambda = -18$

too simple

$\ln \lambda = 0$

$$ln\lambda = -18$$
$$N = 10 \text{ (limited)}$$
$$M = 9$$

$$ln\lambda = 0$$
$$N = 10$$
$$M = 9$$

Q: which one would you choose A, B, C?