



Institute for Advanced Studies
in Basic Sciences
Gava Zang, Zanjan, Iran

Regression with Linear Models

By Dr. Parvin Razzaghi (p.razzaghi@iasbs.ac.ir)

Winter 2024





Slides credit

- Professor:
 - Dr. Parvin Razzaghi
- Original Slides:
 - Amir-massoud Farahman
 - Eric Eaton
- Link:
 - <https://B2n.ir/n25709>
- University:
 - University of Toronto and Vector Institute
- Acknowledgment:
 - Credit for slides goes to many members of the ML Group at the U of T, and beyond, including (recent past): Roger Grosse, Murat Erdogdu, Richard Zemel, Juan Felipe Carrasquilla, Emad Andrews.
 - For any comments or suggestions, please contact: p.razzaghi@iasbs.ac.ir

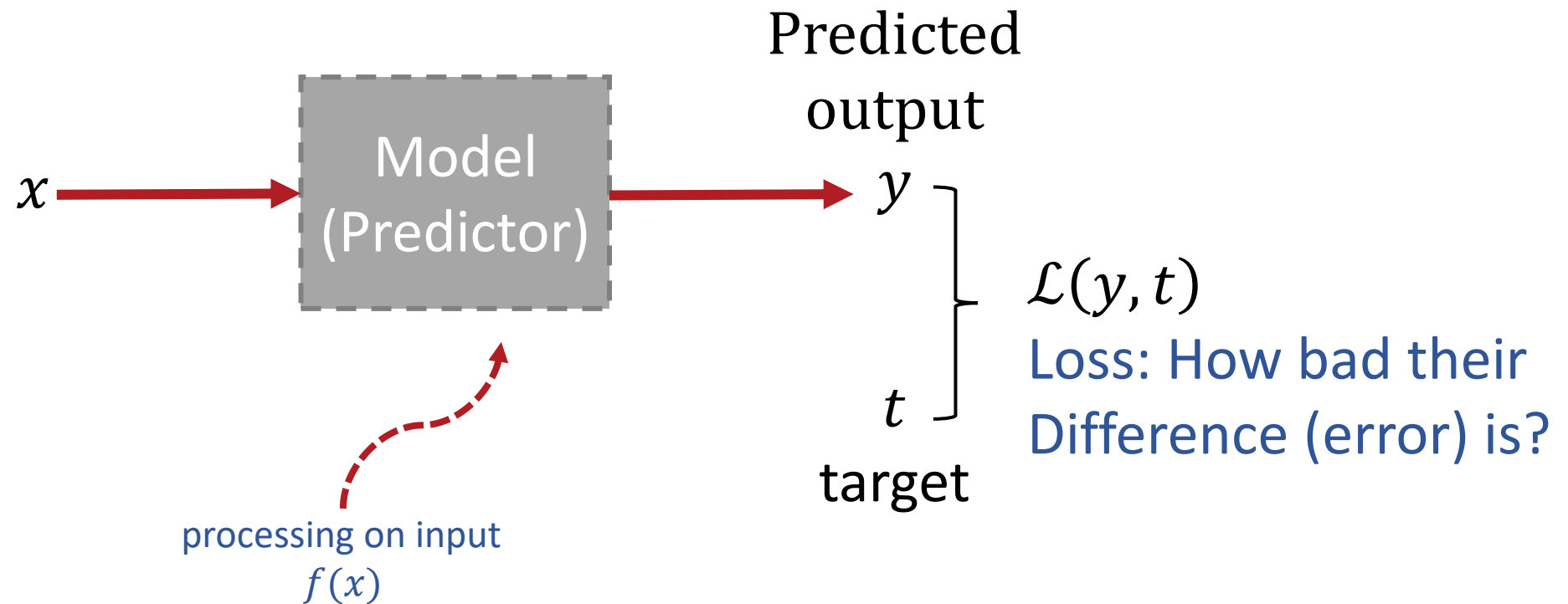


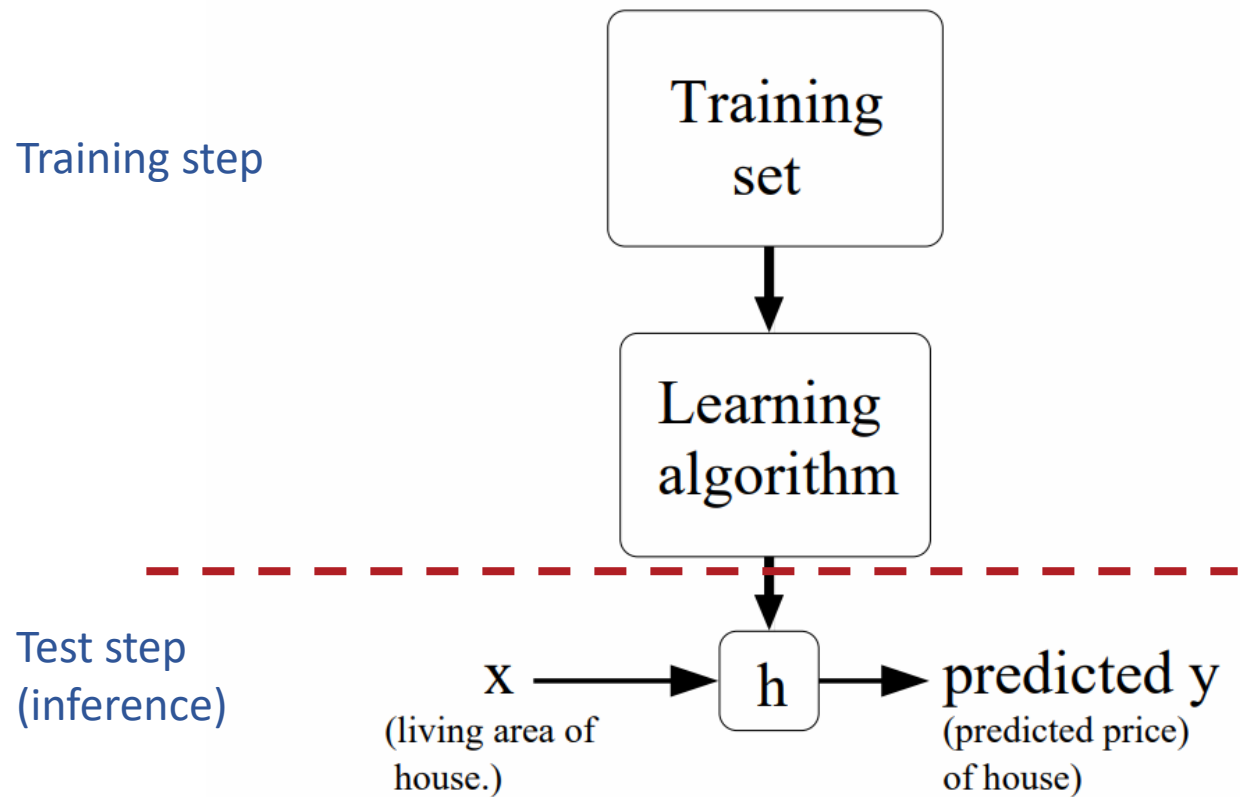
Table of Contents

- Modular Approach to ML
- Regression
 - Linear Regression
 - Basic Expansion
 - Regularization
 - Probabilistic Interpretation of the Squared Error
 - Bias-variance tradeoff
 - Bayesian linear regression



Modular Approach ML Algorithm Design







Modular Approach ML Algorithm Design

- we will take a more **modular** approach:
 - define the problem
 - provide a training dataset
 - choose a **model** describing the relationships between variables of interest
 - define a **loss function** quantifying how bad the fit to the data is
 - (possibly) choose a **regularizer** saying how much we prefer different candidate models (or explanations of data), **before (prior to)** seeing the data
 - fit the model that minimizes the loss function and satisfy the constraint/penalty imposed by the regularizer, possibly using an **optimization algorithm**
- Mixing and matching these modular components gives us a lot of new ML methods.



Define the problem & provide a dataset



The Supervised Learning Setup

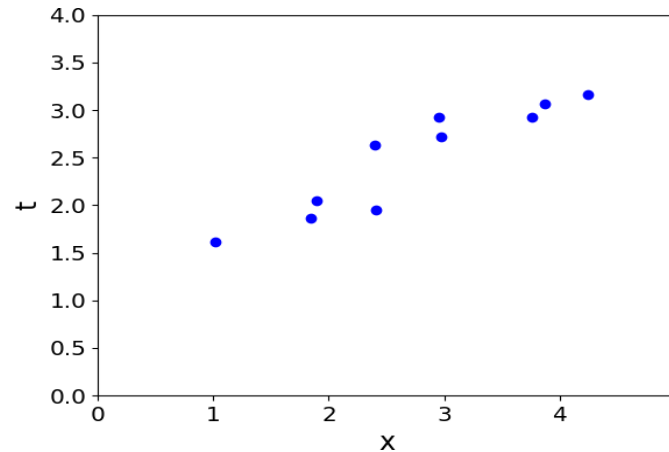
- Features
 - Living area
 - Number of bedrooms
 - ...
- Target variable
 - Price



Living area (feet ²)	#bedrooms	Price (1000\$s)
2104	3	400
1600	3	330
2400	3	369
1416	2	232
3000	4	540
⋮	⋮	⋮



The Supervised Learning Setup



Recall that in supervised learning:

- There is a target $t \in T$ (also called response, outcome, output, class)
- There are features $\mathbf{x} \in X$ (also called inputs or covariates)

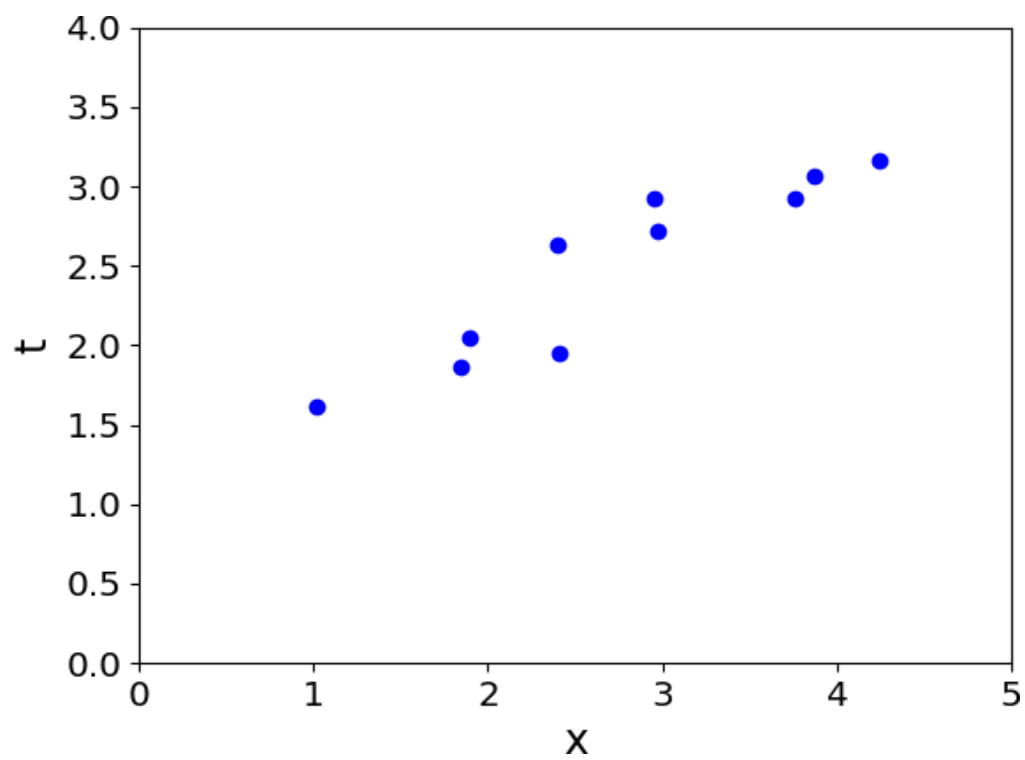
The goal is to learn a function $f : X \rightarrow T$ such that

$$t \approx y = f(\mathbf{x}),$$

based on given data $D = \{(\mathbf{x}^{(i)}, t^{(i)}) \text{ for } i = 1, 2, \dots, N\}$.

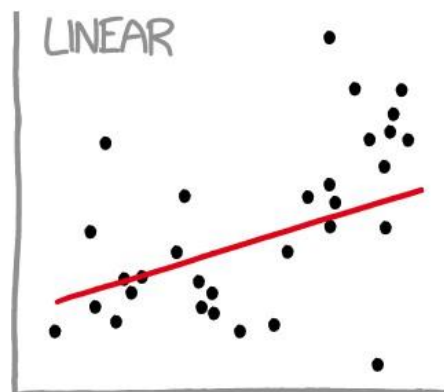


choose a **model** describing the relationships between variables of interest?





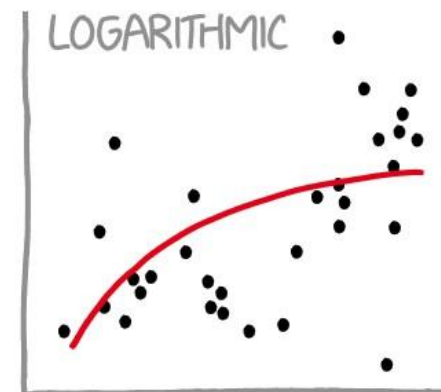
Regression with Linear Models



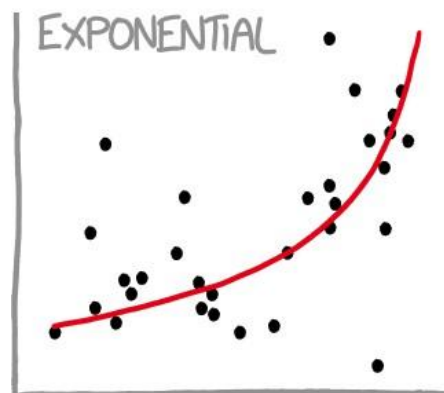
"HEY, I DID A
REGRESSION."



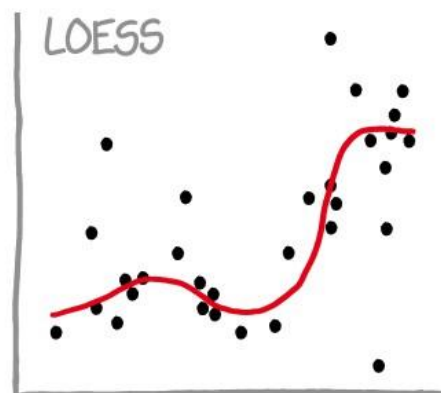
"I WANTED A CURVED
LINE, SO I MADE ONE
WITH MATH."



"LOOK, IT'S
TAPERING OFF!"



"LOOK, IT'S GROWING
UNCONTROLLABLY!"



"I'M SOPHISTICATED, NOT
LIKE THOSE BUMBLING
POLYNOMIAL PEOPLE."



"I'M MAKING A
SCATTER PLOT BUT
I DON'T WANT TO."



Linear Regression – Model

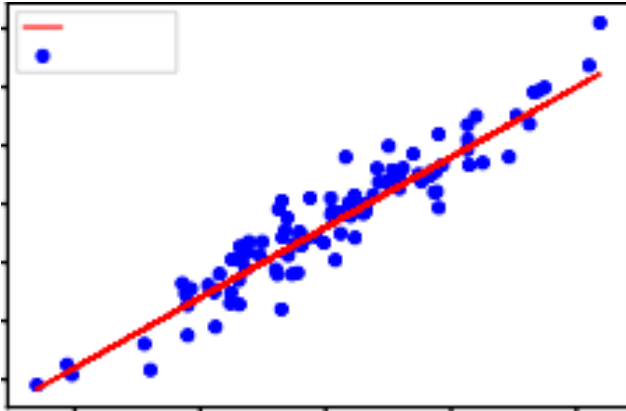
- **Model:** In linear regression, we use linear functions of the inputs $\mathbf{x} = (x_1, \dots, x_D)$ to make predictions y of the target value t :

$$y = f(\mathbf{x}) = \sum_j w_j x_j + b$$

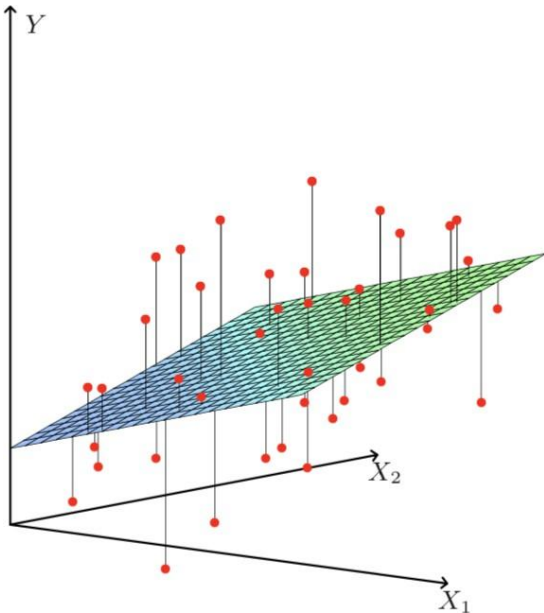
- y is the **prediction**
- \mathbf{w} is the **weights**
- b is the **bias** (or **intercept**) (do not confuse with the bias-variance tradeoff in the next lecture)
- \mathbf{w} and b together are the **parameters**
- We hope that our prediction is close to the target: $y \approx t$.



What is Linear? 1 Feature vs. D Features



- If we have only 1 feature:
 $y = wx + b$ where $w, x, b \in \mathbb{R}$
- y is linear in x

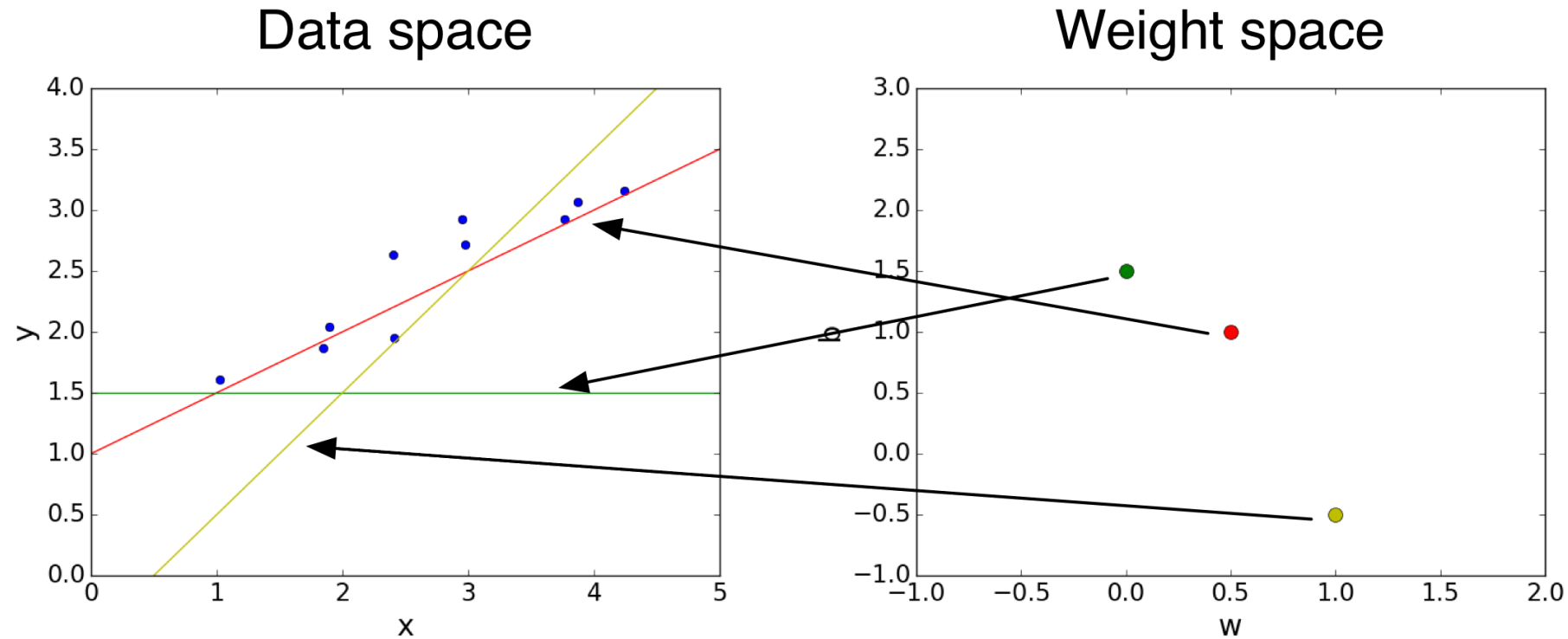


- If we have D features:
 $y = \mathbf{w}^T \mathbf{x} + b$ where $\mathbf{w}, \mathbf{x} \in \mathbb{R}^D, b \in \mathbb{R}$
- y is linear in \mathbf{x}

Relation between the prediction y and inputs \mathbf{x} is linear in both cases.



Weight Space vs. Data Space



Recall that:

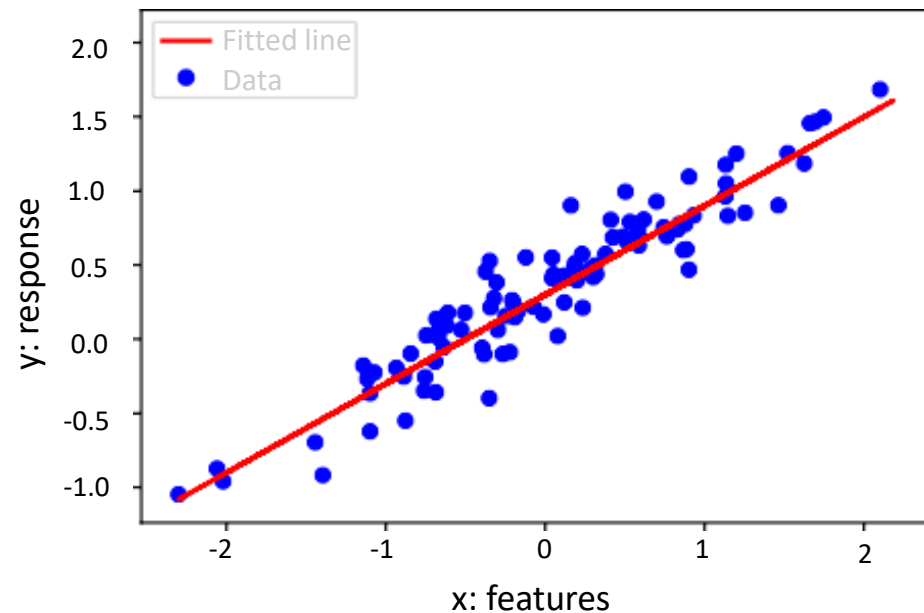
$$y = f(x) = \sum_j w_j x_j + b$$



Linear Regression

We have a dataset $D = \{(x^{(i)}, t^{(i)})\}_{i=1}^N$ where:

- $x^{(i)} = (x_1^{(i)}, x_2^{(i)}, \dots, x_D^{(i)})^T \in \mathbb{R}^D$ are the inputs, e.g., age, height
- $t^{(i)} \in \mathbb{R}$ is the target or response, e.g., income
- predict $t^{(i)}$ with a linear function of $x^{(i)}$:



- $t^{(i)} \approx y^{(i)} = \mathbf{w}^T \mathbf{x}^{(i)} + b$
- Find the “best” line (\mathbf{w}, b) .
- **Q:** How should we find the **best** line?

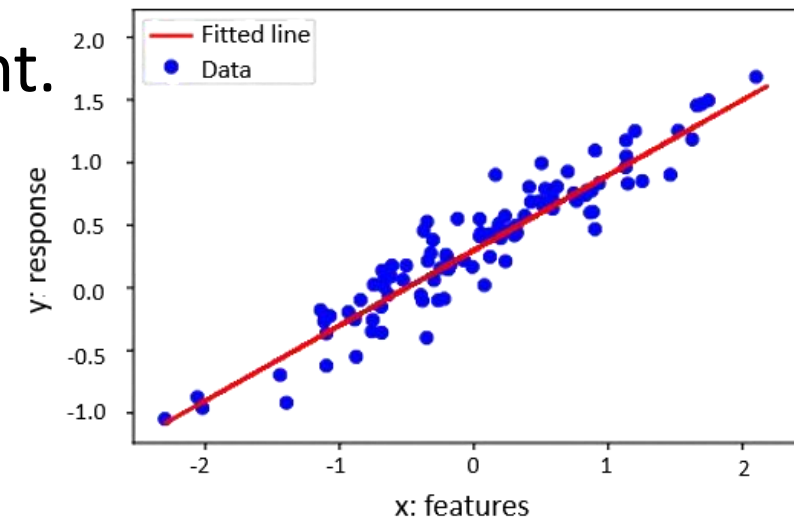


Linear Regression – Loss Function

- How to quantify the quality of the fit to data?
- A **loss function** $\mathcal{L}(y, t)$ defines how bad it is if, for some input x , the algorithm predicts y , but the target is actually t .
- **Squared error loss function:**

$$\mathcal{L}(y, t) = \frac{1}{2} (y - t)^2$$

- $y - t$ is the **residual**, and we want to make its magnitude small
- The $\frac{1}{2}$ factor is just to make the calculations convenient.





Linear Regression – Loss Function

- **Cost function:** loss function averaged over all training examples

$$\mathcal{J}(\mathbf{w}, b) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(y^{(i)}, t^{(i)})$$

$$= \frac{1}{2N} \sum_{i=1}^N (y^{(i)} - t^{(i)})^2$$

$$= \frac{1}{2N} \sum_{i=1}^N (\mathbf{w}^T x^{(i)} + b - t^{(i)})^2$$



Linear Regression – Loss Function

- To find the best fit, we find a model (parameterized by its weights \mathbf{w} and b) that minimizes the cost:

$$\underset{(\mathbf{w}, b)}{\text{minimize}} \mathcal{J}(\mathbf{w}, b) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(y^{(i)}, t^{(i)})$$

- The terminology is not universal. Some might call “loss” **pointwise loss** and the “cost function” the **empirical loss** or **average loss**.



Vector Notation

- We can organize all the training examples into a **design matrix** \mathbf{X} with one row per training example, and all the targets into the **target vector** \mathbf{t} .

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}^{(1)T} \\ \mathbf{x}^{(2)T} \\ \mathbf{x}^{(3)T} \end{bmatrix} = \begin{bmatrix} 8 & 0 & 3 & 0 \\ 6 & -1 & 5 & 3 \\ 2 & 5 & -2 & 8 \end{bmatrix}$$

One feature across
all training examples

One training
example (vector)

- Computing the predictions for the whole dataset:

$$\mathbf{X}\mathbf{w} + b\mathbf{1} = \begin{bmatrix} \mathbf{w}^T \mathbf{x}^{(1)} + b \\ \vdots \\ \mathbf{w}^T \mathbf{x}^{(N)} + b \end{bmatrix} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{bmatrix} = \mathbf{y}$$



Vectorization

- Computing the squared error cost across the whole dataset:

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

$$\mathcal{J} = \frac{1}{2N} \|\mathbf{y} - \mathbf{t}\|^2$$

- Note that sometimes we may use $\mathcal{J} = \frac{1}{2N} \|\mathbf{y} - \mathbf{t}\|^2$, without $\frac{1}{N}$ normalizer. That would correspond to the sum of losses, and not the average loss. That does not matter as the minimizer does not depend on N .



Vectorization

- We can also add a column of 1s to the design matrix, combine the bias and the weights, and conveniently write

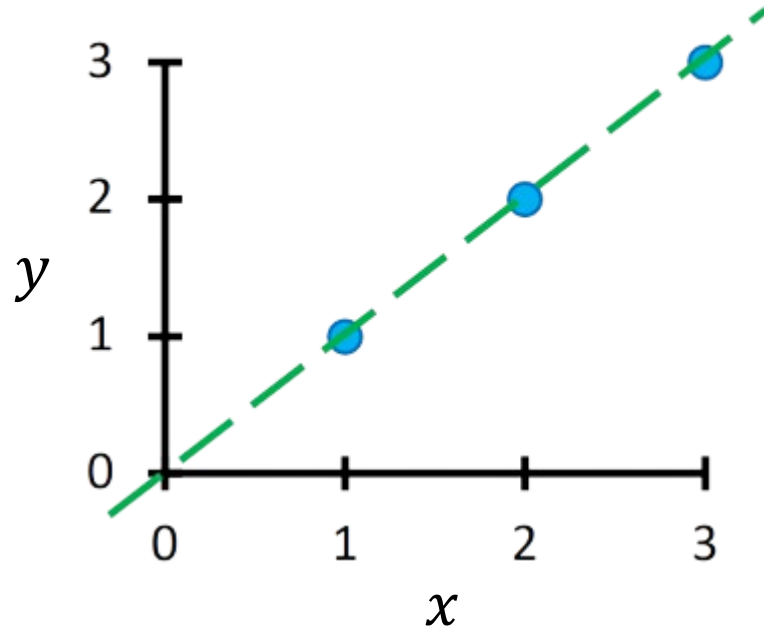
$$X = \begin{bmatrix} 1 & [x^{(1)}]^T \\ 1 & [x^{(2)}]^T \\ 1 & \vdots \end{bmatrix} \in \mathbb{R}^{N \times D+1} \text{ and } w = \begin{bmatrix} b \\ w_1 \\ w_2 \\ \vdots \end{bmatrix} \in \mathbb{R}^{D+1}$$

Then, our predictions reduce to $\mathbf{y} = \mathbf{X}\mathbf{w}$.



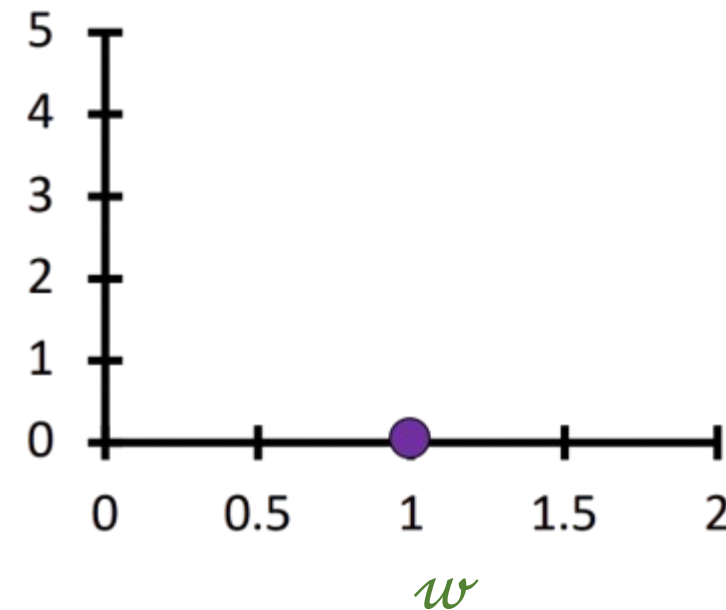
Intuition on Minimizing MSE Loss

- Consider $x \in \mathbb{R}$ and $w \in \mathbb{R}$



$$w = 1$$

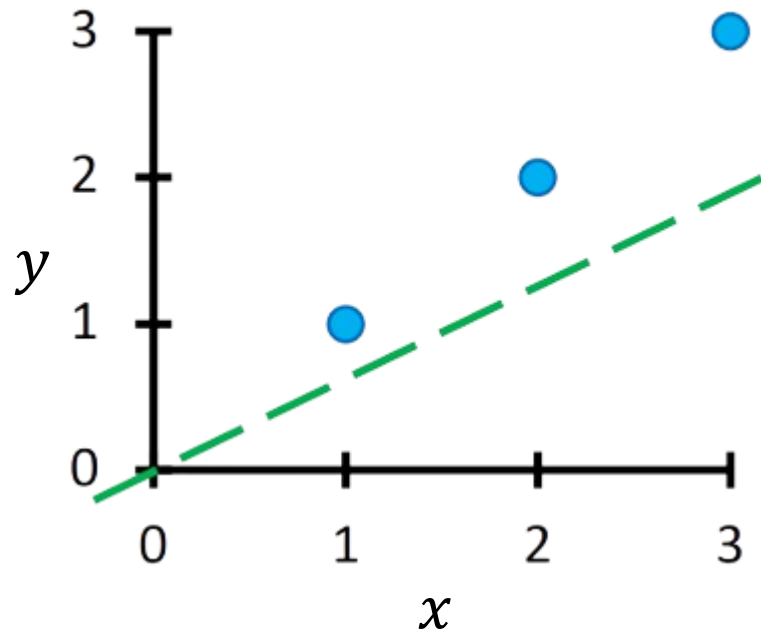
$$J(w, D)$$



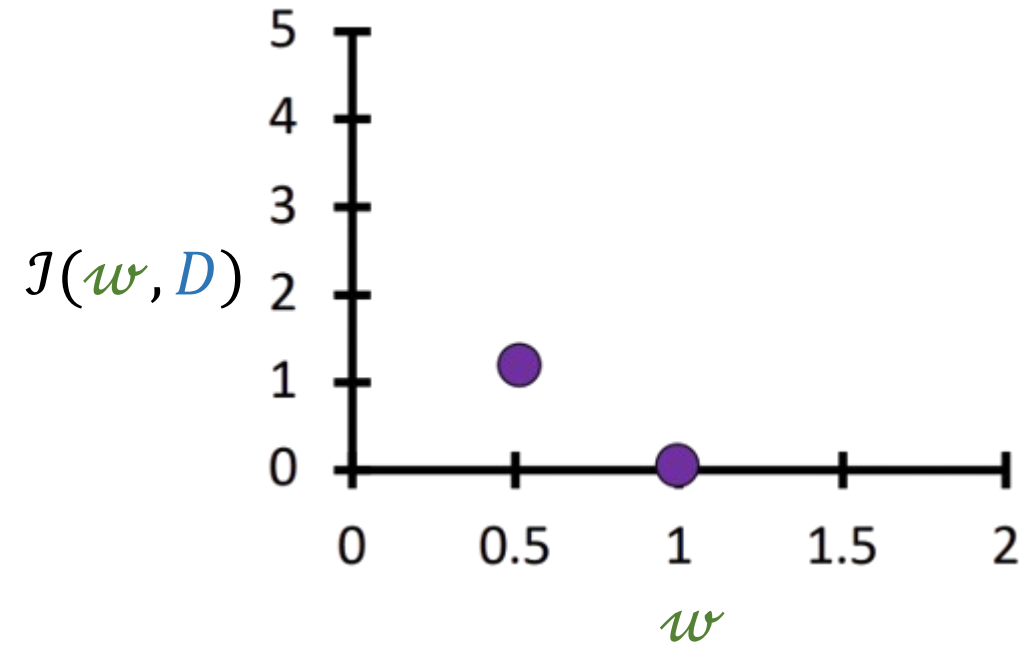


Intuition on Minimizing MSE Loss

- Consider $x \in \mathbb{R}$ and $w \in \mathbb{R}$



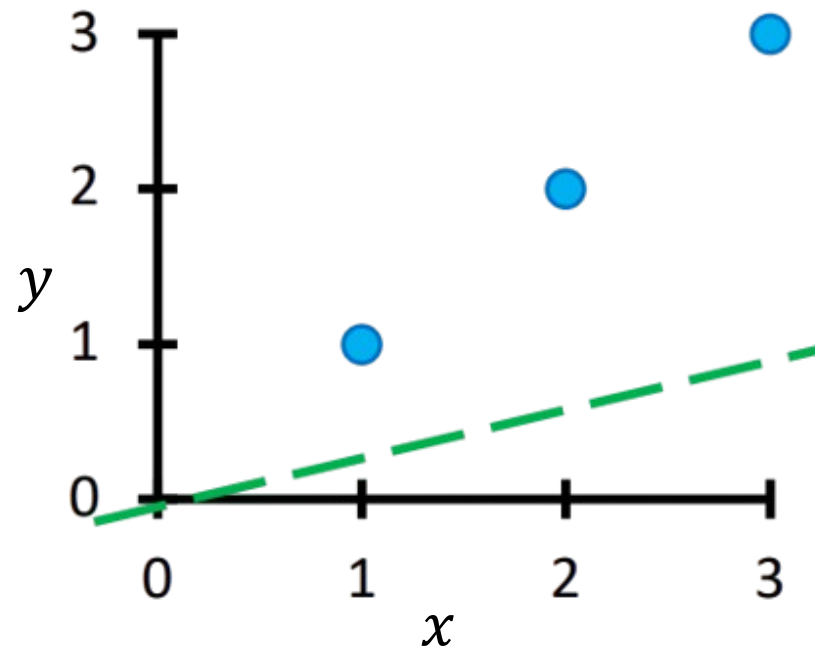
$$w = 0.5$$





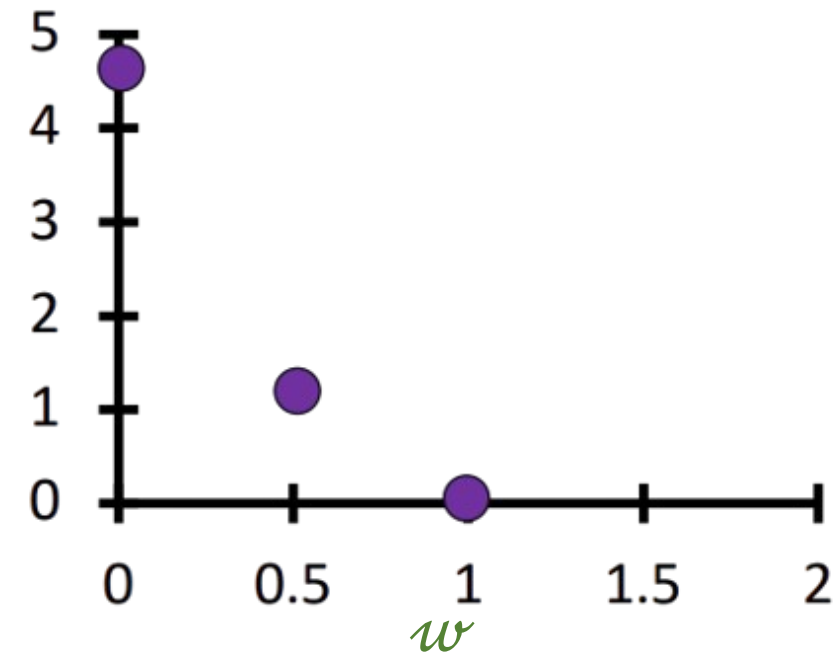
Intuition on Minimizing MSE Loss

- Consider $x \in \mathbb{R}$ and $w \in \mathbb{R}$



$$w = 0.25$$

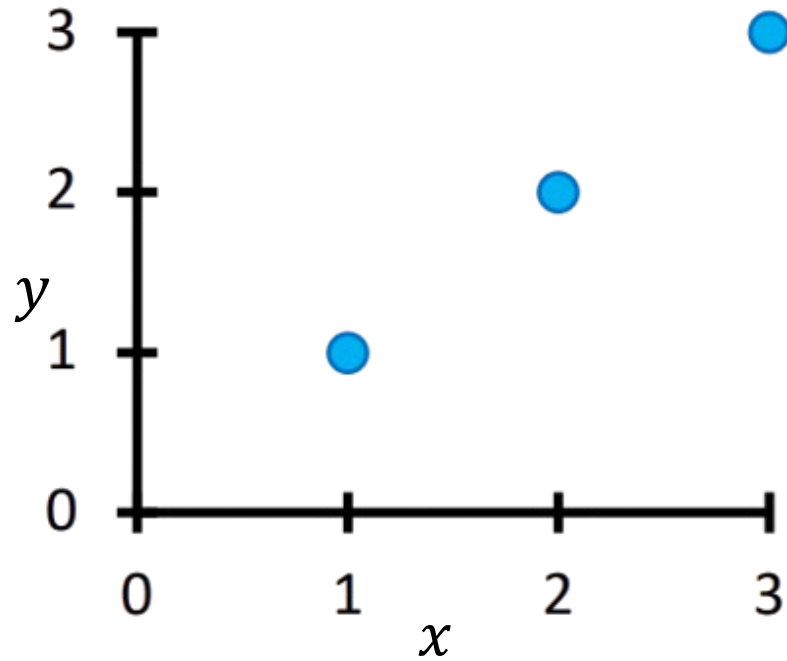
$$J(w, D)$$



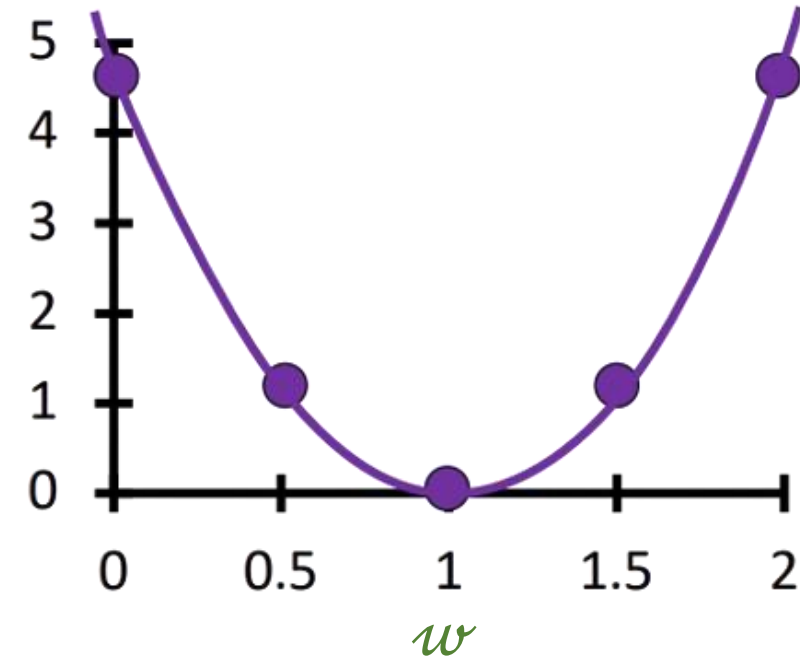


Intuition on Minimizing MSE Loss

- Consider $x \in \mathbb{R}$ and $w \in \mathbb{R}$



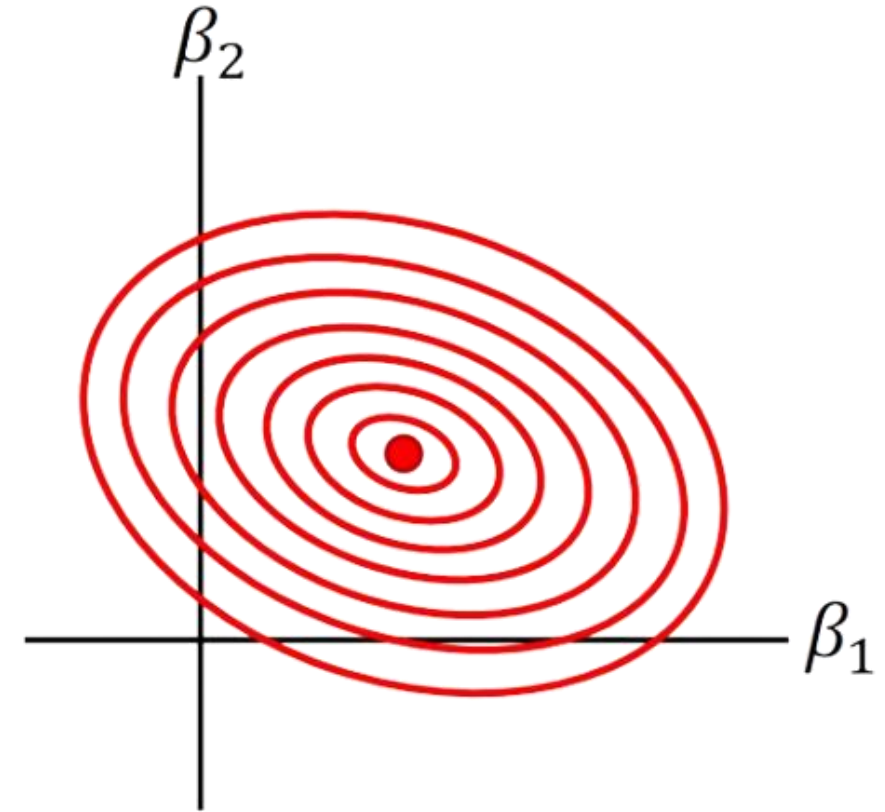
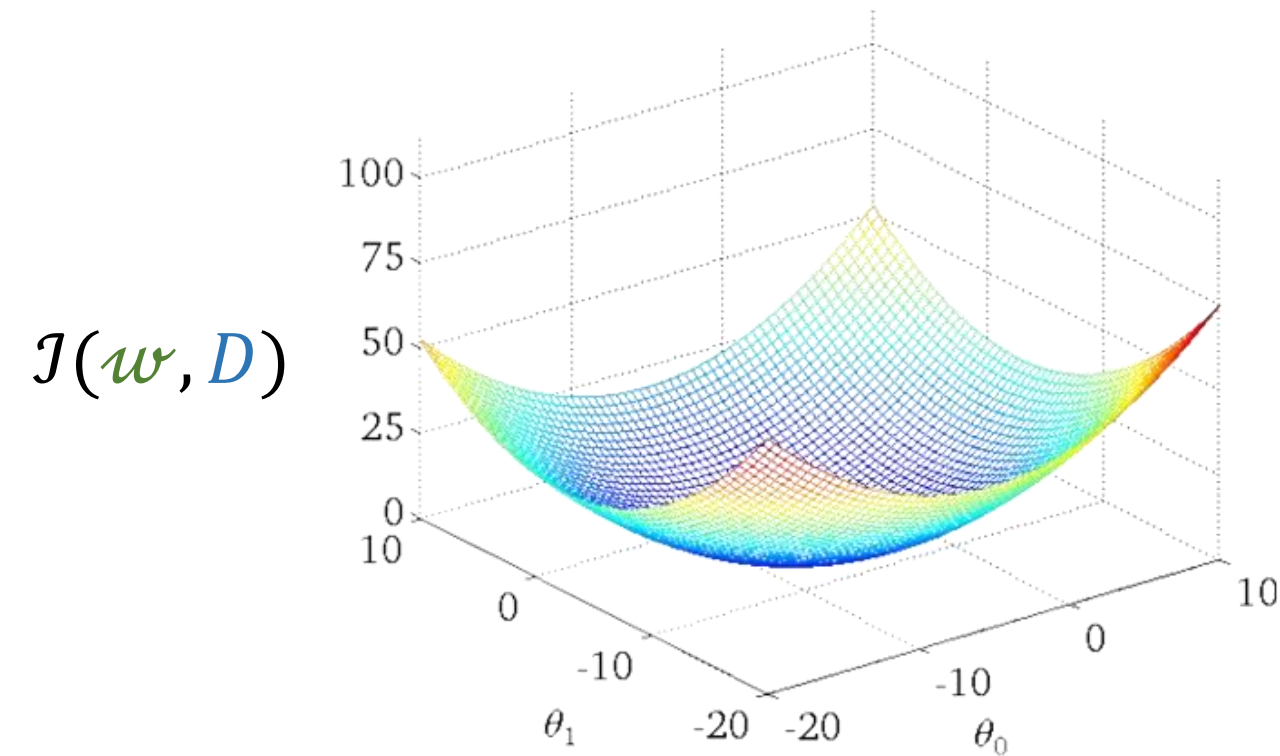
$J(w, D)$





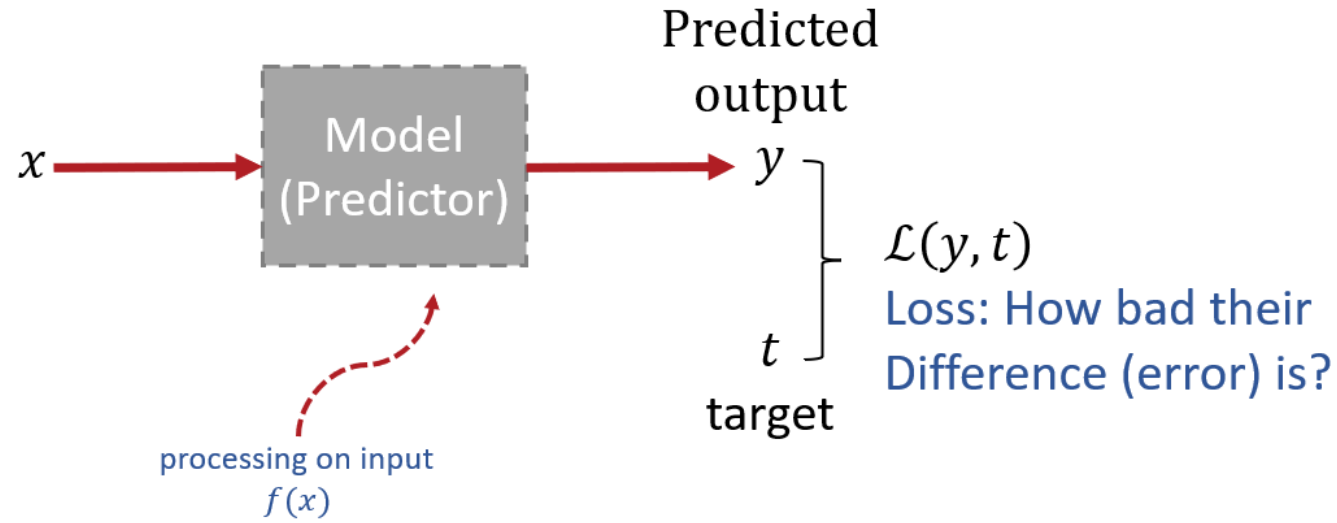
Intuition on Minimizing MSE Loss

- **Convex** (“bowl shaped”) in general





Solving the Minimization Problem



- We defined a model (linear).
- We defined a loss and the cost function to be minimized.
- **Q:** How should we solve this minimization problem?



Solving the Minimization Problem

- Recall from your calculus class: minimum of a differentiable function (if it exists) occurs at a **critical point**, i.e., point where the derivative is zero.
- Multivariate generalization: set the partial derivatives to zero (or equivalently the gradient).
- We would like to find a point where the gradient is (close to) zero. How can we do it?
- Sometimes it is possible to directly find the parameters that make the gradient zero in a closed-form. We call this the **direct solution**.
- We may also use optimization techniques that iteratively get us closer to the solution. We will get back to this soon.



Direct Solution

- **Partial derivatives:** derivatives of a multivariate function with respect to (w.r.t.) one of its arguments.

$$\frac{\partial}{\partial x_1} f(x_1, x_2) = \lim_{h \rightarrow 0} \frac{f(x_1 + h, x_2) - f(x_1, x_2)}{h}$$



Direct Solution

- To compute, take the single variable derivatives, pretending the other arguments are constant.
- Example: partial derivatives of the prediction y with respect to weight w_j and bias b :

$$\frac{\partial y}{\partial w_j} = \frac{\partial}{\partial w_j} \left[\sum_{j'} w_{j'} x_{j'} + b \right] = x_j$$

$$\frac{\partial y}{\partial b} = \frac{\partial}{\partial b} \left[\sum_{j'} w_{j'} x_{j'} + b \right] = 1$$



Direct Solution

- **The derivative of loss:** We apply the chain rule: first we take the derivative of the loss L w.r.t. output y of the model, and then the derivative of the output y w.r.t. a parameter of the model such as w_j or b :

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial w_j} &= \frac{d\mathcal{L}}{dy} \frac{\partial y}{\partial w_j} \\ &= \frac{d}{dy} \left[\frac{1}{2} (y - t)^2 \right] \cdot x_j \\ &= (y - t) x_j \\ \frac{\partial \mathcal{L}}{\partial b} &= y - t\end{aligned}$$



Direct Solution

- Cost derivatives (average over data points):

$$\frac{\partial \mathcal{J}}{\partial w_j} = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - t^{(i)}) x_j^{(i)}$$

$$\frac{\partial \mathcal{J}}{\partial b} = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - t^{(i)})$$



Direct Solution

- Recall that the output y is a function of the parameters as $y = \mathbf{w}^T \mathbf{x}$.
- The minimum of the cost function must occur at a point where the partial derivatives are zero, i.e.,

$$\nabla_{\mathbf{w}} \mathcal{J} = 0 \iff \frac{\partial \mathcal{J}}{\partial w_j} = 0 \quad (\forall j), \quad \frac{\partial \mathcal{J}}{\partial b} = 0.$$

- If $\partial \mathcal{J} / \partial w_j \neq 0$, you could reduce the cost by changing w_j .



Direct Solution

If we follow this recipe, we get that we have to set the gradient of

$\mathcal{J} = \frac{1}{2N} \|\mathbf{y} - \mathbf{t}\|^2$, with $\mathbf{y} = \mathbf{X}\mathbf{w}$ (bias absorbed in \mathbf{X}) equal to zero.

We have

$$\mathcal{J} = \frac{1}{2N} (\mathbf{X}\mathbf{w} - \mathbf{t})^T (\mathbf{X}\mathbf{w} - \mathbf{t}),$$

So

$$\nabla_{\mathbf{w}} \mathcal{J} = \frac{1}{N} \mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{t}) = 0 \implies (\mathbf{X}^T \mathbf{X})\mathbf{w} = \mathbf{X}^T \mathbf{t}.$$

This is a **linear system of equations**.

Q: What are the dimensions of each component?



Direct Solution

- Assuming that $\mathbf{X}^T \mathbf{X}$ is invertible, the optimal weights are

$$\mathbf{w}^{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}.$$

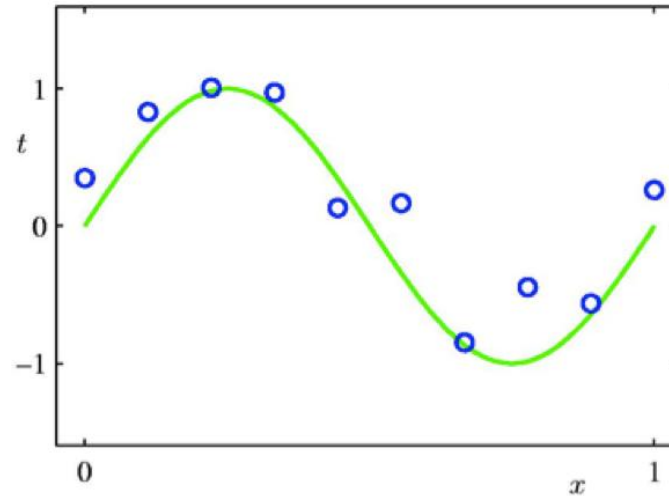
This solution is also called **Ordinary Least Squares (OLS)** solution.

At an arbitrary point \mathbf{x} , our prediction is $y = \mathbf{w}^{LS^T} \mathbf{x}$.

Q: What happens if $\mathbf{X}^T \mathbf{X}$ is not invertible?

Basis Expansion (Feature Mapping)

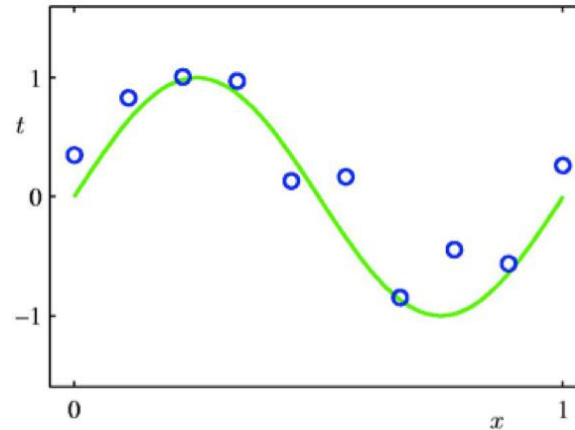
- The relation between the input and output may not be linear.



- We can still use linear regression by mapping the input feature to another space using **basis expansion** (or **feature mapping**) $\psi(\mathbf{x}) : \mathbb{R}^D \rightarrow \mathbb{R}^d$ and treat the mapped feature (in \mathbb{R}^d) as the input of a linear regression procedure.
- Let us see how it works when $x \in \mathbb{R}$ and we use polynomial feature mapping.



Polynomial Feature Mapping



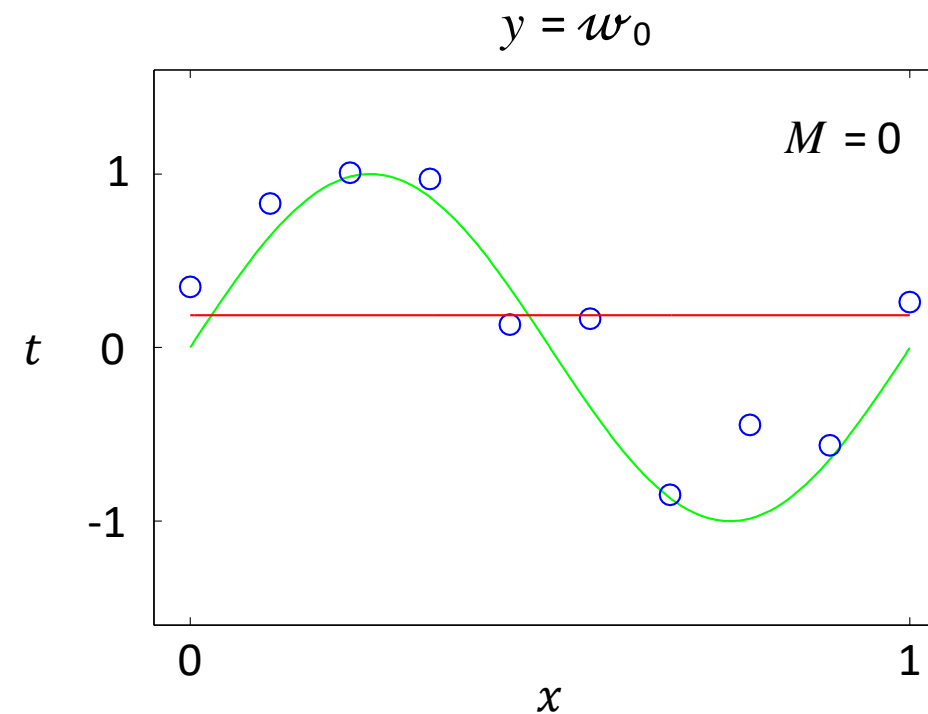
- Fit the data using a degree- M polynomial function of the form:

$$y = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{i=0}^M w_i x^i$$

- The feature mapping is $\psi(x) = [1, x, x^2, \dots, x^M]^T$.
- We can still use the linear regression framework with least squares loss to find \mathbf{w} since $y = \psi(x)^T \mathbf{w}$ is linear in w_0, w_1, \dots .
- In general, ψ can be any function. Another example: Fourier map
$$\psi = [1, \sin(2\pi x), \cos(2\pi x), \sin(4\pi x), \cos(4\pi x), \sin(6\pi x), \cos(6\pi x), \dots]^T.$$

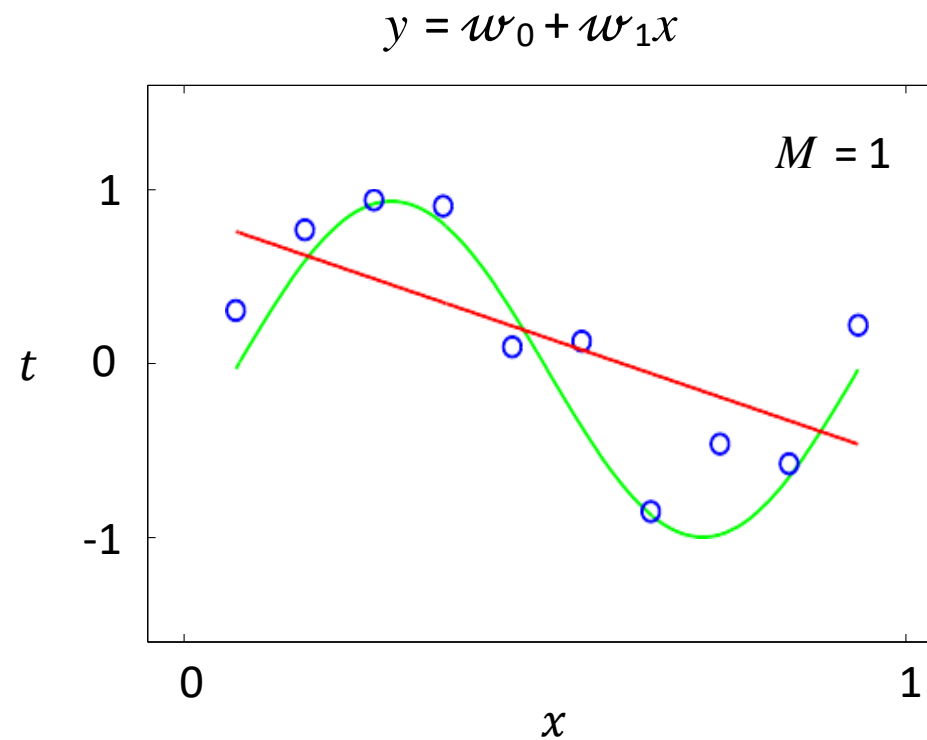


Polynomial Feature Mapping with $M = 0$



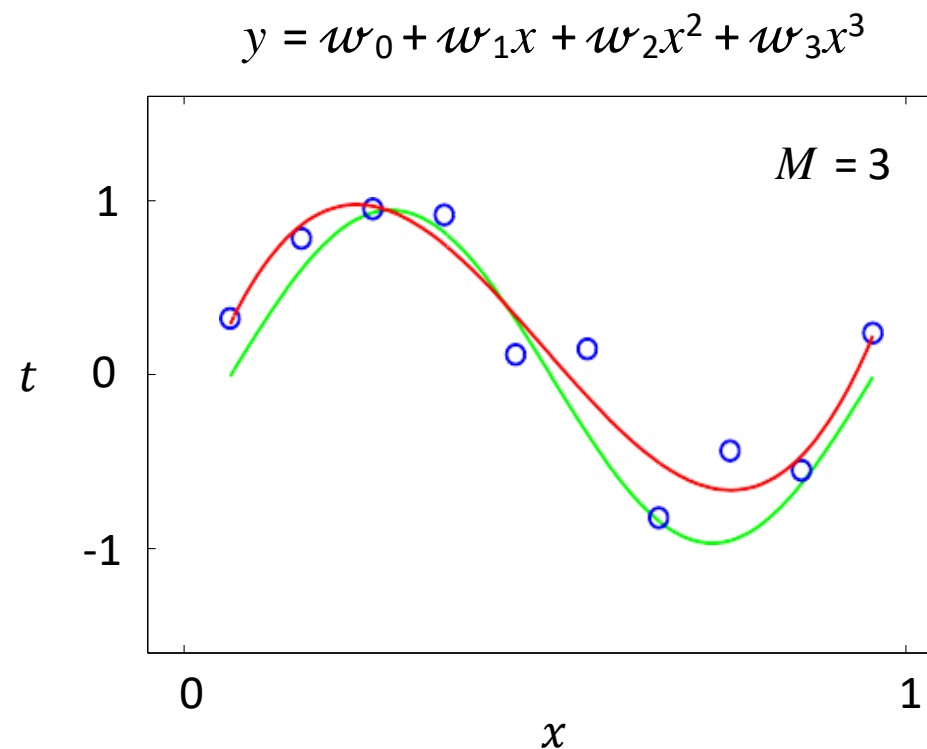


Polynomial Feature Mapping with $M = 1$



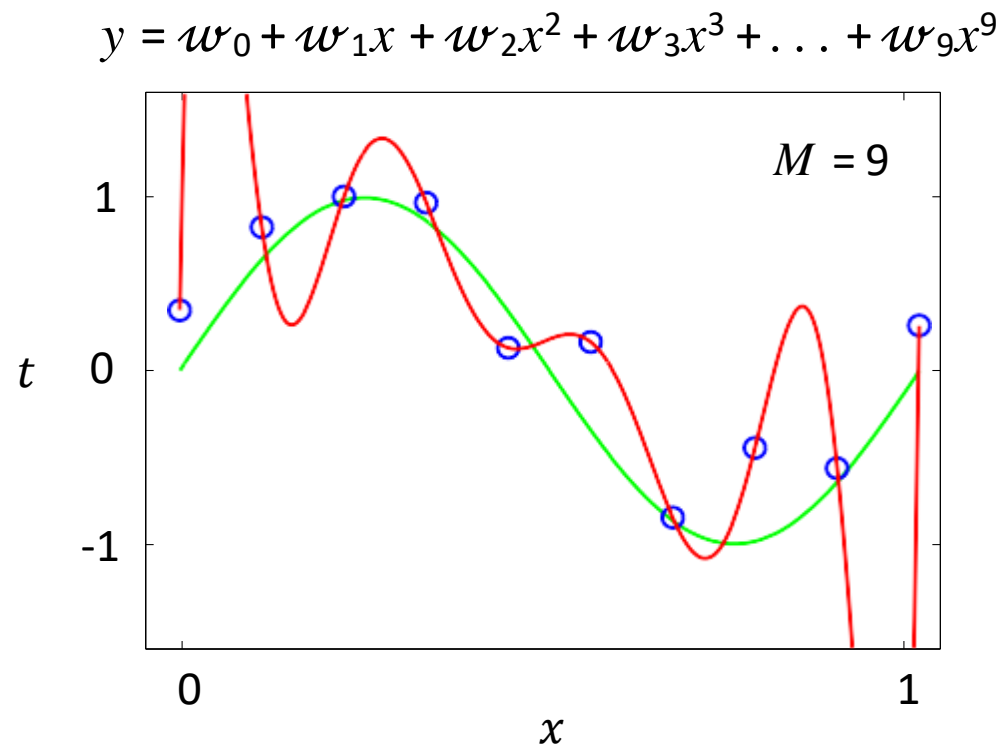


Polynomial Feature Mapping with $M = 3$



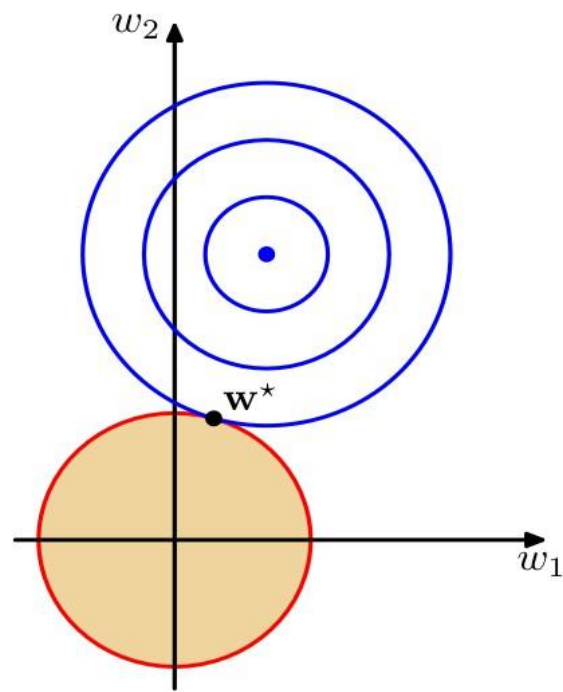


Polynomial Feature Mapping with $M = 9$



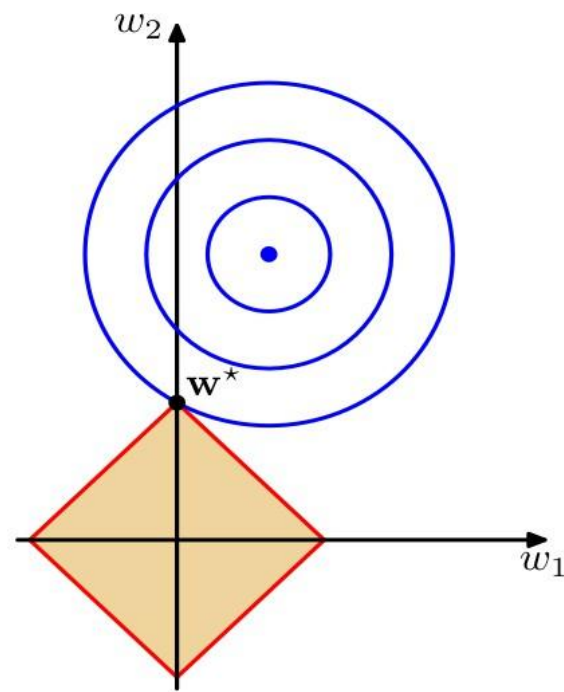


Model Complexity and Regularization



L2 regularization

$$\mathcal{R} = \sum_i w_i^2$$



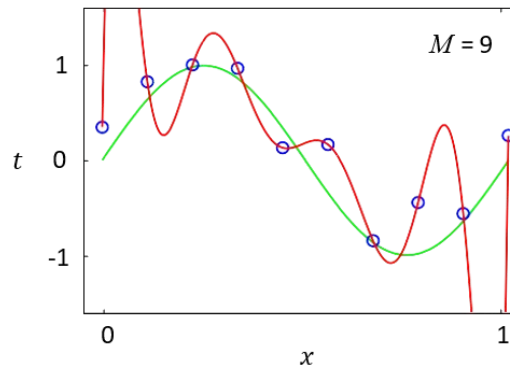
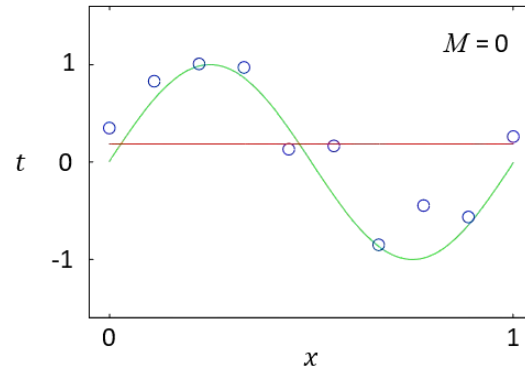
L1 regularization

$$\mathcal{R} = \sum_i |w_i|$$

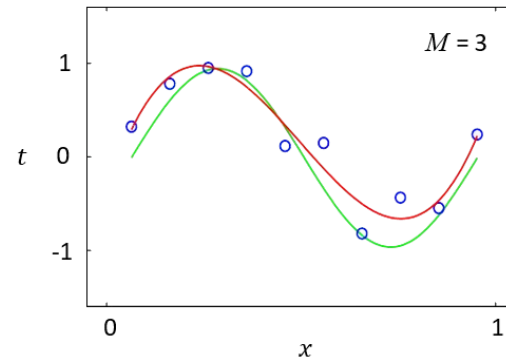


Model Complexity and Generalization

- **Underfitting** ($M=0$): model is too simple — does not fit the data.
- **Overfitting** ($M=9$): model is too complex — fits perfectly.



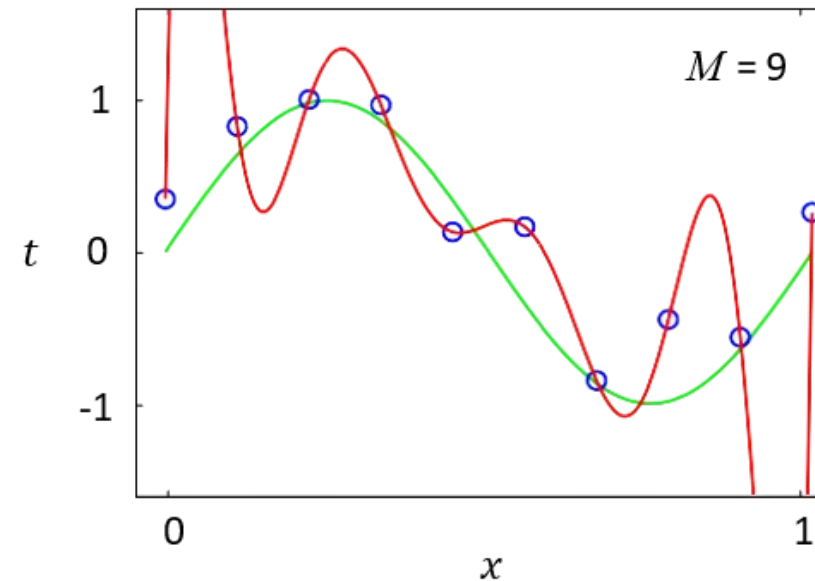
- **Good model** ($M=3$): Achieves small test error (generalizes well).





Model Complexity and Generalization

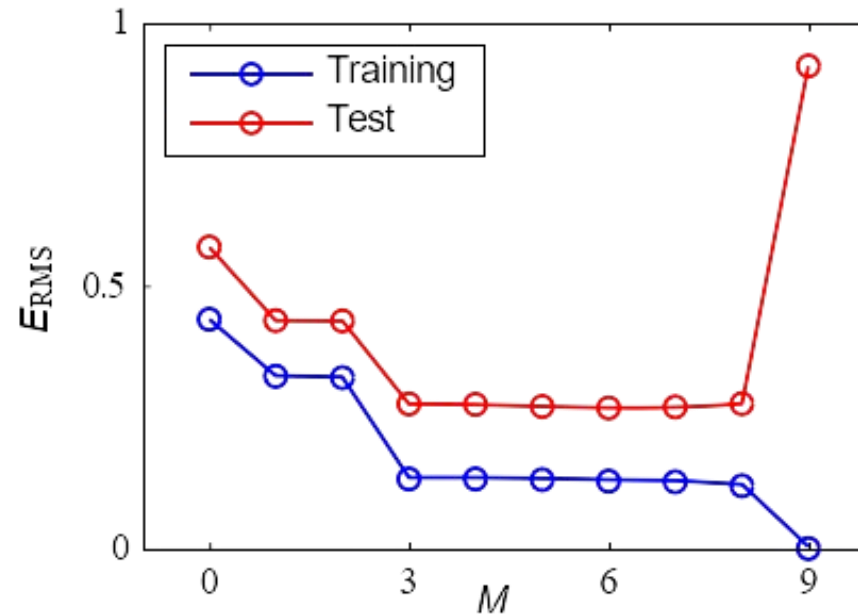
	$M = 0$	$M = 1$	$M = 3$	$M = 9$
w_0^*	0.19	0.82	0.31	0.35
w_1^*		-1.27	7.99	232.37
w_2^*			-25.43	-5321.83
w_3^*			17.37	48568.31
w_4^*				-231639.30
w_5^*				640042.26
w_6^*				-1061800.52
w_7^*				1042400.18
w_8^*				-557682.99
w_9^*				125201.43



- As M increases, the magnitude of coefficients gets larger.
- For $M = 9$, the coefficients have become finely tuned to the data.
- Between data points, the function exhibits large oscillations.



Model Complexity and Generalization



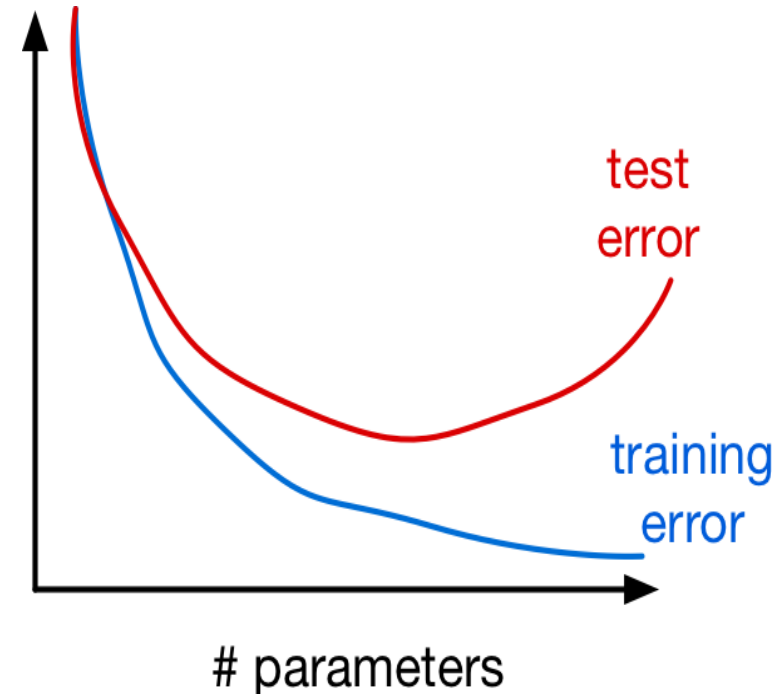
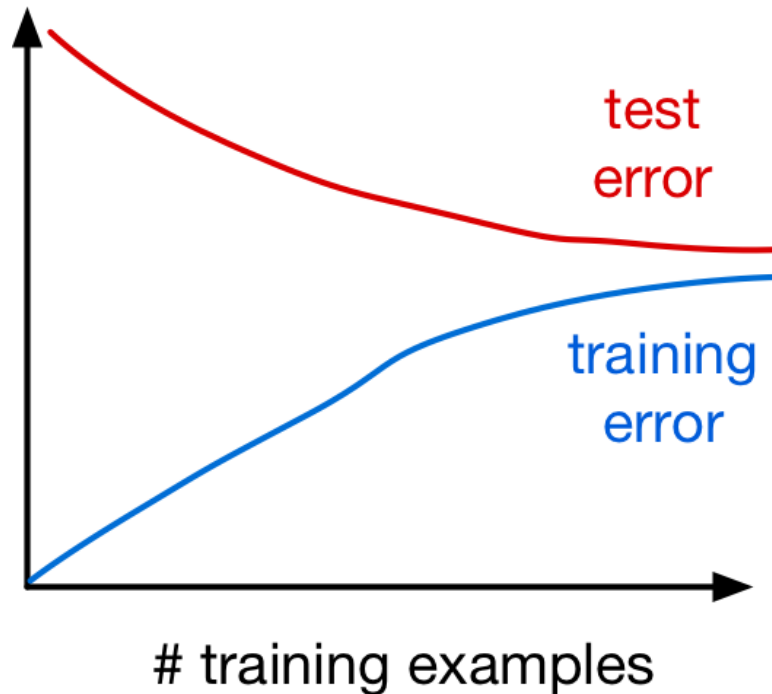
As the degree M of the polynomial increases

- the training errors decreases;
- the test error, however, initially decreases, but then increases.



Model Complexity and Generalization

- Training and test error as a function of # training examples and parameters:





Regularization for Controlling the Model Complexity

- The degree of the polynomial M controls the **complexity** of the model.
- The value of M is a hyperparameter for polynomial expansion, just like K in KNN or the depth of a tree in a decision tree. We can tune it using a validation set.
- Restricting the number of parameters of a model (M here) is a crude approach to control the complexity of the model.
- A better solution: keep the number of parameters of the model large, but enforce “**simpler**” solutions within the same space of parameters.
- This is done through **regularization** or **penalization**.
 - **Regularizer** (or **penalty**): a function that quantifies how much we prefer one hypothesis vs. another, prior to seeing the data.
- **Q:** How?!



ℓ_2 or (L^2) Regularization

- We can encourage the weights to be small by choosing the ℓ_2 (or L^2) of the weights as our regularizer or **penalty**:

$$\mathcal{R} = \frac{1}{2} \|w\|_2^2 = \frac{1}{2} \sum_j w_j^2.$$

- Note: To be precise, we are regularizing the *squared* l_2 norm.
- The regularized cost function makes a tradeoff between fit to the data and the norm of the weights:

$$\mathcal{J}_{reg}(w) = \mathcal{J}(w) + \lambda \mathcal{R}(w) = \mathcal{J}(w) + \frac{\lambda}{2} \sum_j w_j^2.$$



ℓ_2 or (L^2) Regularization

- The regularized cost function

$$\mathcal{J}_{reg}(w) = \mathcal{J}(w) + \lambda \mathcal{R}(w) = \mathcal{J}(w) + \frac{\lambda}{2} \sum_j w_j^2.$$

- The basic idea is that “simpler” functions have weights w with smaller ℓ_2 -norm and we prefer them to functions with larger ℓ_2 -norm.
 - **Intuition:** Large weights makes the function f have more abrupt changes as a function of the input x ; it will be less smooth.
- If you fit training data poorly, \mathcal{J} is large. If the fitted weights have high values, \mathcal{R} is large.
- Large λ penalizes weight values more.
- Here λ is a hyperparameter that we can tune with a validation set.



ℓ_2 Regularized Least Squares: Ridge Regression

For the least squares problem, we have $\mathcal{J}(w) = \frac{1}{2N} \|Xw - t\|^2$.

- When $\lambda > 0$ (with regularization), regularized cost gives

$$\begin{aligned} w_{\lambda}^{Ridge} &= \underset{w}{\operatorname{argmin}} \mathcal{J}_{reg}(w) = \underset{w}{\operatorname{argmin}} \frac{1}{2N} \|Xw - t\|_2^2 + \frac{\lambda}{2} \|w\|_2^2 \\ &= (X^T X + \lambda N \mathbf{I})^{-1} X^T t. \end{aligned}$$

- The case of $\lambda = 0$ (no regularization) reduces to the least squares solution!
- **Q:** What happens when $\lambda \rightarrow \infty$?
- Note that it is also common to formulate this problem as

$\underset{w}{\operatorname{argmin}} \|Xw - t\|_2^2 + \frac{\lambda}{2} \|w\|_2^2$ in which case the solution is

$$w_{\lambda}^{Ridge} = (X^T X + \lambda \mathbf{I})^{-1} X^T t.$$



Lasso and the ℓ_1 Regularization

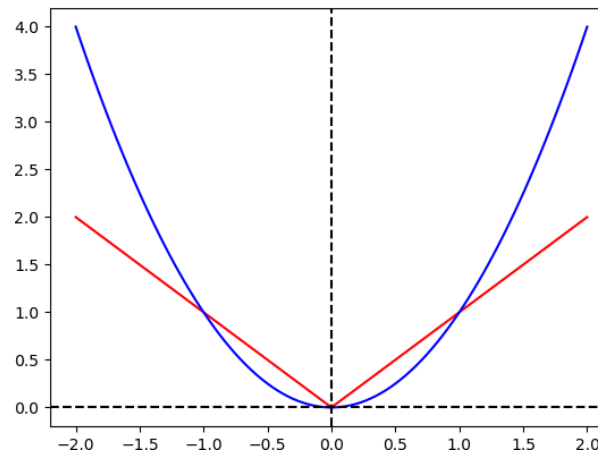
- The ℓ_1 norm, or sum of absolute values, is another regularizer:

$$\mathcal{R}(w) = \|w\|_1 = \sum_j |w_j|$$

- The **Lasso** (Least Absolute Shrinkage and Selection Operator) is

$$\min_w \|X\mathbf{w} - \mathbf{t}\|_2^2 + \lambda \|\mathbf{w}\|_1.$$

- It can be shown that Lasso encourages weights to be exactly zero.
 - **Q:** When is this helpful?





Ridge vs. Lasso – Geometric Viewpoint

- We presented regularization as a penalty on the weights, in which we solve

$$\min_w \mathcal{J}(w) + \lambda \mathcal{R}(w)$$

- We can also write an equivalent form as a constraint optimization:

$$\operatorname{argmin}_w \mathcal{J}(w)$$

$$s.t. \quad \mathcal{R}(w) \leq \mu$$

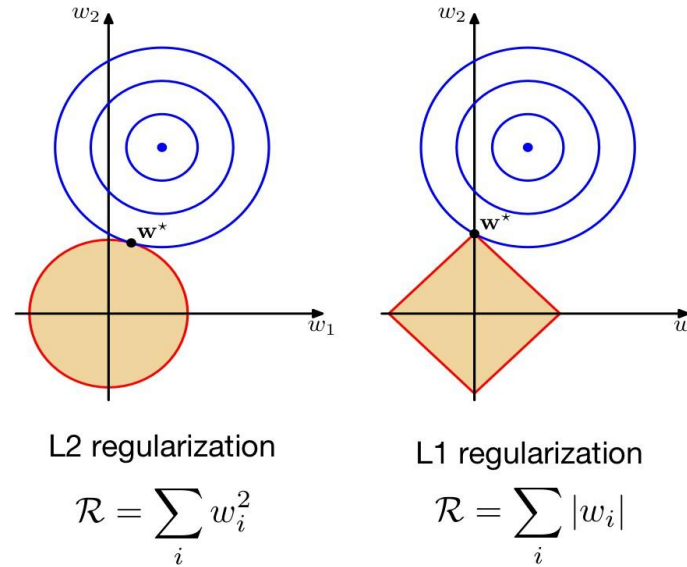
for a corresponding value of μ

- The Ridge regression and the Lasso can then be written as

$$\operatorname{argmin}_w \|X\mathbf{w} - \mathbf{t}\|_2^2$$

$$s.t. \quad \|\mathbf{w}\|_p \leq \mu \quad (Lasso: p = 1; Ridge: p = 2)$$

Ridge vs. Lasso – Geometric Viewpoint



- The set $\{\mathbf{w}: \|X\mathbf{w} - \mathbf{t}\|_2^2 \leq \varepsilon\}$ defines ellipsoids of ε cost in the weights space.
- The set $\{\mathbf{w}: \|\mathbf{w}\|_p \leq \mu\}$ defines the constraint on weights defined by the regularizer.
- The solution would be the smallest ε for which these two sets intersect.
- For $p = 1$, the diamond-shaped constraint set has corners. When the intersection happens at a corner, some of the weights are zero.
- For $p = 2$, the disk-shaped constraint set does not have corners. It does not induce any zero weights.



Alternative loss function

- Mean absolute error

$$\frac{1}{n} \sum_{i=1}^n |y_i - t_i|$$

- Mean relative error

$$\frac{1}{n} \sum_{i=1}^n \frac{|y_i - t_i|}{|t_i|}$$

- R^2 score:

- “Coefficient of determination”
- Higher is better, $R^2 = 1$ is perfect

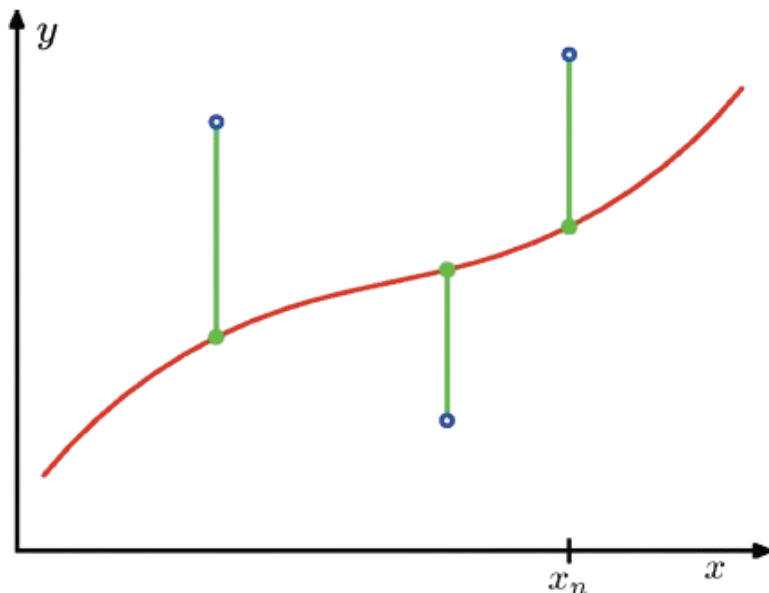
$$1 - \frac{MSE}{Variance}$$



Probabilistic Interpretation of the Squared Error

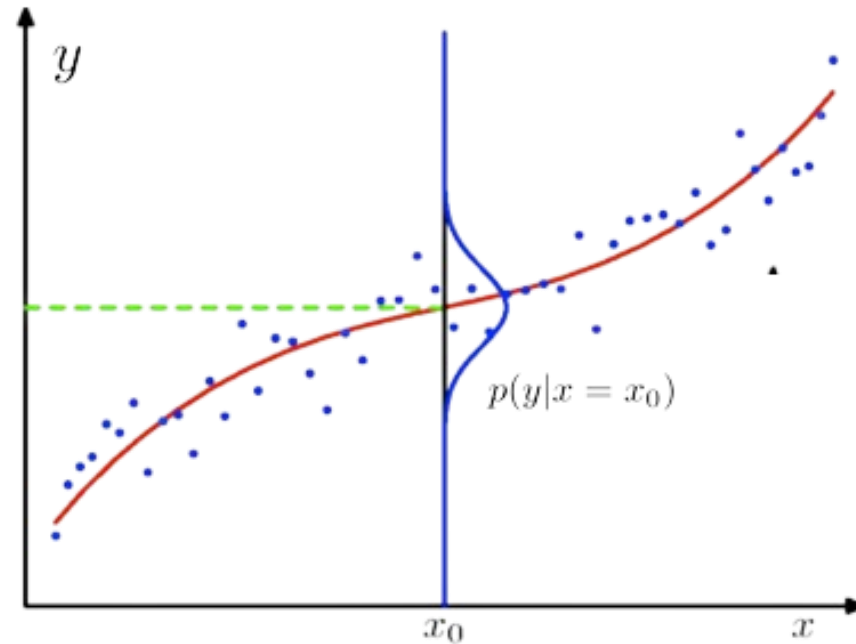
- For the least squares: we minimize the sum of the squares of the errors between the predictions for each data point $\mathbf{x}^{(i)}$ and the corresponding target values $t^{(i)}$, i.e.,

$$\underset{(\mathbf{w}, \mathbf{w}_0)}{\text{minimize}} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}^{(i)} + b - t^{(i)})^2$$



- $t \approx \mathbf{x}^T \mathbf{w} + b, (\mathbf{w}, b) \in \mathbb{R}^D \times \mathbb{R}$
- We measure the quality of the fit using the squared error loss. Why?
- Even though the squared error loss is intuitive, we did not justify it.
- We provide a probabilistic perspective here.
- There are other justifications too; we get to them in the Bias-Variance decomposition lecture.

Probabilistic Interpretation of the Squared Error



- Suppose that our model arose from a statistical model ($b = 0$ for simplicity):

$$t^{(i)} = \mathbf{w}^T \mathbf{x}^{(i)} + \varepsilon^{(i)}$$

where $\varepsilon^{(i)} \sim \mathcal{N}(0, \sigma^2)$ is independent of the input $\mathbf{x}^{(i)}$.

- Thus, $t^{(i)} | \mathbf{x}^{(i)} \sim p(t | \mathbf{x}^{(i)}, \mathbf{w}) = \mathcal{N}(\mathbf{w}^T \mathbf{x}^{(i)}, \sigma^2)$.



Probabilistic Interpretation of the Squared Error: Maximum Likelihood Estimation

- Suppose that the input data $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$ are given and the outputs are independently drawn from

$$t^{(i)} \sim p(y|\mathbf{x}^{(i)}, \mathbf{w}).$$

- with an unknown parameter \mathbf{w} . So the dataset is

$$D = \{(\mathbf{x}^{(1)}, t^{(1)}), \dots, (\mathbf{x}^{(N)}, t^{(N)})\}.$$



Probabilistic Interpretation of the Squared Error: Maximum Likelihood Estimation

- The **likelihood** function is $\Pr(D|\mathbf{w})$.
- The **maximum likelihood estimation (MLE)** is based on the “**principle**” suggesting that we have to find a parameter $\hat{\mathbf{w}}$ that maximizes the likelihood, i.e.,

$$\hat{\mathbf{w}} \leftarrow \underset{\mathbf{w}}{\operatorname{argmax}} \Pr(D|\mathbf{w}) .$$

Maximum likelihood estimation: after observing the data samples $(\mathbf{x}^{(t)}, t^{(i)})$ for $i = 1, 2, \dots, N$, we should choose \mathbf{w} that maximizes the likelihood.



Probabilistic Interpretation of the Squared Error: Maximum Likelihood Estimation

- For independent samples, the likelihood function of samples D is the product of their likelihoods

$$p(t^{(1)}, t^{(2)}, \dots, t^{(N)} | \mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}, w) = \prod_{i=1}^N p(t^{(i)} | \mathbf{x}^{(i)}, w) = L(w).$$

- Product of N terms is not easy to minimize.
- Taking log reduces it to a sum. Two objectives are equivalent since log is strictly increasing.
- Maximizing the likelihood is equivalent to minimizing the **negative log-likelihood**:



Probabilistic Interpretation of the Squared Error: Maximum Likelihood Estimation

- Maximizing the likelihood is equivalent to minimizing the **negative log-likelihood**:

$$\ell(w) = -\log L(w) = -\log \prod_{i=1}^N p(t^{(i)} | \mathbf{x}^{(i)}; \mathbf{w}) = -\sum_{i=1}^n \log p(t^{(i)} | \mathbf{x}^{(i)}; \mathbf{w}).$$



Probabilistic Interpretation of the Squared Error: Maximum Likelihood Estimation

Maximum Likelihood Estimator (MLE)

After observing $\mathbf{z}^{(i)} = (\mathbf{x}^{(i)}, t^{(i)})$ for $i = 1, \dots, N$ independent and identically distributed (i.i.d.) samples from $p(\mathbf{z}, \mathbf{w})$, MLE is

$$\mathbf{w}^{MLE} = \underset{\mathbf{w}}{\operatorname{argmin}} l(\mathbf{w}) = - \sum_{i=1}^N \log p(t^{(i)} | \mathbf{x}^{(i)}; \mathbf{w}).$$



Probabilistic Interpretation of the Squared Error: From MLE to Squared Error

- Suppose that our model arose from a statistical model:

$$t^{(i)} = \mathbf{w}^T \mathbf{x}^{(i)} + \varepsilon^{(i)}$$

where $\varepsilon^{(i)} \sim \mathcal{N}(0, \sigma^2)$ is independent of anything else.

- $p(t^{(i)} | \mathbf{x}^{(i)}, w) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2} (t^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2\right\}$
- $\log p(t^{(i)} | \mathbf{x}^{(i)}, w) = \frac{1}{2\sigma^2} (t^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2 - \log(\sqrt{2\pi\sigma^2})$



Probabilistic Interpretation of the Squared Error: From MLE to Squared Error

- The MLE solution is

$$\mathbf{w}^{MLE} = \underset{\mathbf{w}}{\operatorname{argmin}} \mathcal{L}(\mathbf{w}) = \frac{1}{2\sigma^2} \sum_{i=1}^N (t^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2 + \mathcal{C}$$

- As \mathcal{C} and σ do not depend on \mathbf{w} , they do not contribute to the minimization.

$\mathbf{w}^{MLE} = \mathbf{w}^{LS}$ when we work with Gaussian densities.



Probabilistic Interpretation of the Squared Error: From MLE to Squared Error

- Suppose that our model arose from a statistical model:

$$t^{(i)} = \mathbf{w}^T \mathbf{x}^{(i)} + \varepsilon^{(i)}$$

where $\varepsilon^{(i)}$ comes from the **Laplace** distribution, that is, the distribution of $\varepsilon^{(i)}$ has density.

$$\frac{1}{2b} \exp \left(\frac{|t^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)}|}{b} \right)$$



Probabilistic Interpretation of the Squared Error: From MLE to Squared Error

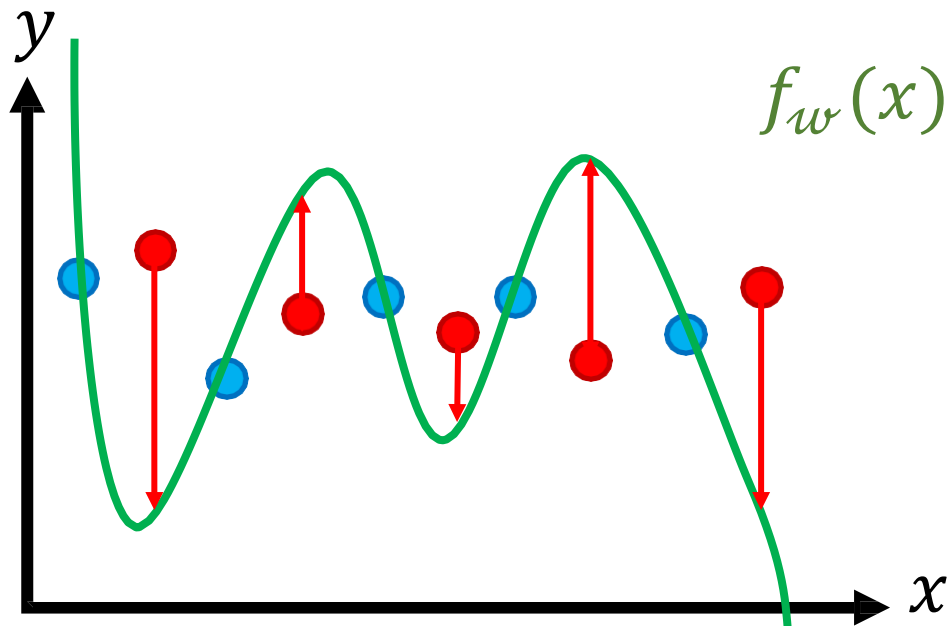
- **Q:** What is the loss in the MLE?
 - Choice 1: $\frac{1}{N} \sum_{i=1}^N |t^{(i)} - w^T x^{(i)}|^{1/2}$
 - Choice 2: $\frac{1}{N} \sum_{i=1}^N (t^{(i)} - w^T x^{(i)})$
 - Choice 3: $\frac{1}{N} \sum_{i=1}^N |t^{(i)} - w^T x^{(i)}|$
 - Choice 4: $\frac{1}{N} \left| \sum_{i=1}^N t^{(i)} - w^T x^{(i)} \right|$
- **Q:** Can you think of an application area with non-Gaussian probabilistic model?



Bias-Variance Tradeoff

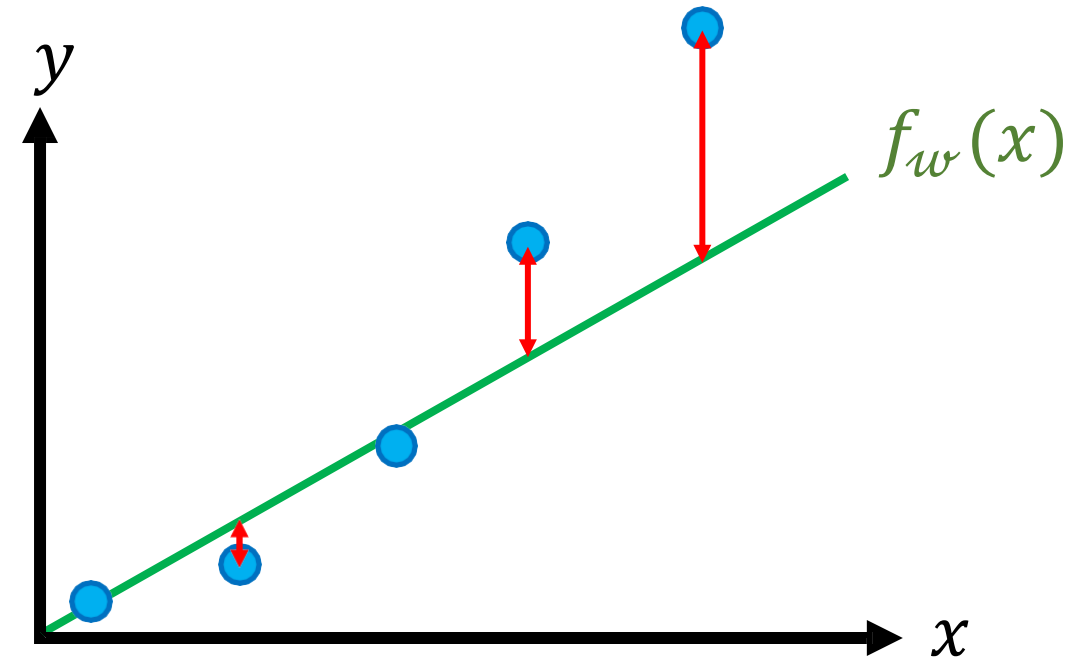
- **Overfitting (high variance)**

- High capacity model capable of fitting complex data
- Insufficient data to constrain it



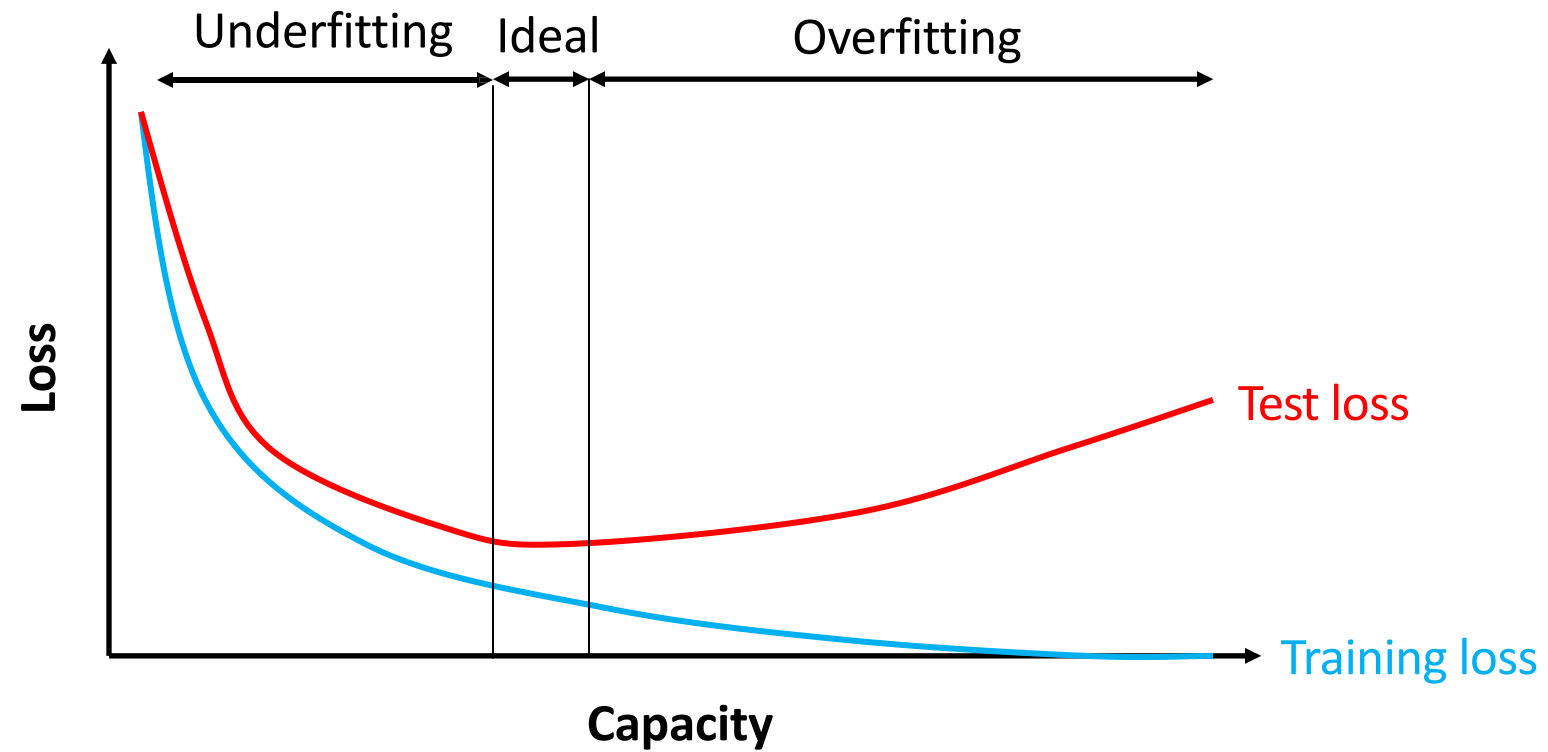
- **Underfitting (high bias)**

- Low capacity model that can only fit simple data
- Sufficient data but poor fit





Bias-Variance Tradeoff





The Bias-Variance Decomposition₍₁₎

- Recall the *expected squared loss*,

$$\mathbb{E}[L] = \int \{y(x) - h(x)\}^2 p(x) dx + \underbrace{\iint \{h(x) - t\}^2 P(x, t) dx dt}_{\text{noise}}$$

where

$$h(x) = \mathbb{E}[t|x] = \int t p(t|x) dt.$$

The second term of $E[L]$ corresponds to the noise inherent in the random variable t .

Q: What about the first term?



The Bias-Variance Decomposition₍₂₎

- Suppose we were given multiple data sets, each of size N . Any particular data set, \mathcal{D} , will give a particular function $y(x; \mathcal{D})$. We then have

$$\begin{aligned} & \{y(x; \mathcal{D}) - h(x)\}^2 \\ &= \{y(x; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(x; \mathcal{D})] + \mathbb{E}_{\mathcal{D}}[y(x; \mathcal{D})] - h(x)\}^2 \\ &= \{y(x; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(x; \mathcal{D})]\}^2 + \{\mathbb{E}_{\mathcal{D}}[y(x; \mathcal{D})] - h(x)\}^2 \\ &+ 2\{y(x; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(x; \mathcal{D})]\} \{\mathbb{E}_{\mathcal{D}}[y(x; \mathcal{D})] - h(x)\} \end{aligned}$$



The Bias-Variance Decomposition₍₃₎

- Taking the expectation over D yields

$$\begin{aligned} & \mathbb{E}_{\mathcal{D}}[\{y(x; \mathcal{D}) - h(x)\}^2] \\ &= \underbrace{\{\mathbb{E}_{\mathcal{D}}[y(x; \mathcal{D})] - h(x)\}^2}_{(\text{bias})^2} + \underbrace{\mathbb{E}_{\mathcal{D}}[\{y(x; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(x; \mathcal{D})]\}^2]}_{\text{variance}} \end{aligned}$$



The Bias-Variance Decomposition₍₄₎

- Thus we can write

$$\textit{expected loss} = (\textit{bias})^2 + \textit{variance} + \textit{noise}$$

where

$$(\textit{bias})^2 = \int \{\mathbb{E}_{\mathcal{D}}[y(x; \mathcal{D})] - h(x)\}^2 p(x) dx$$

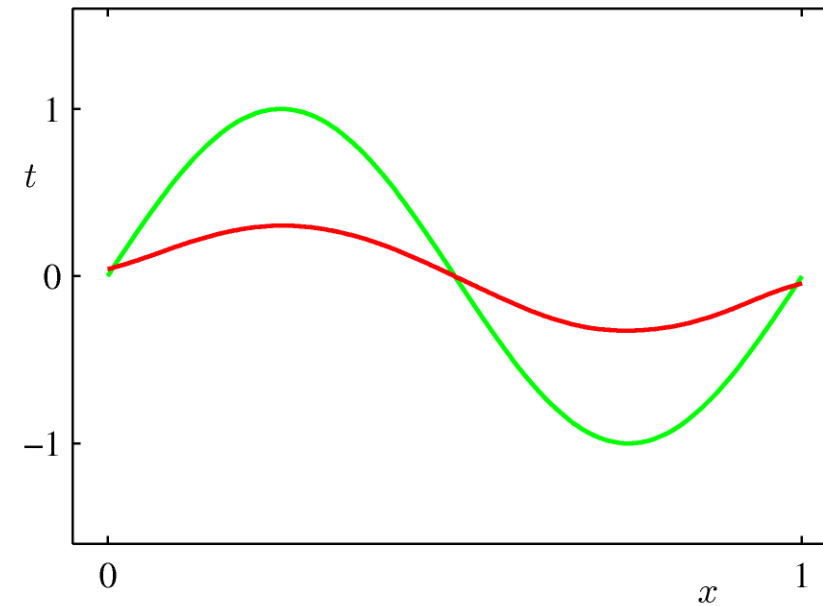
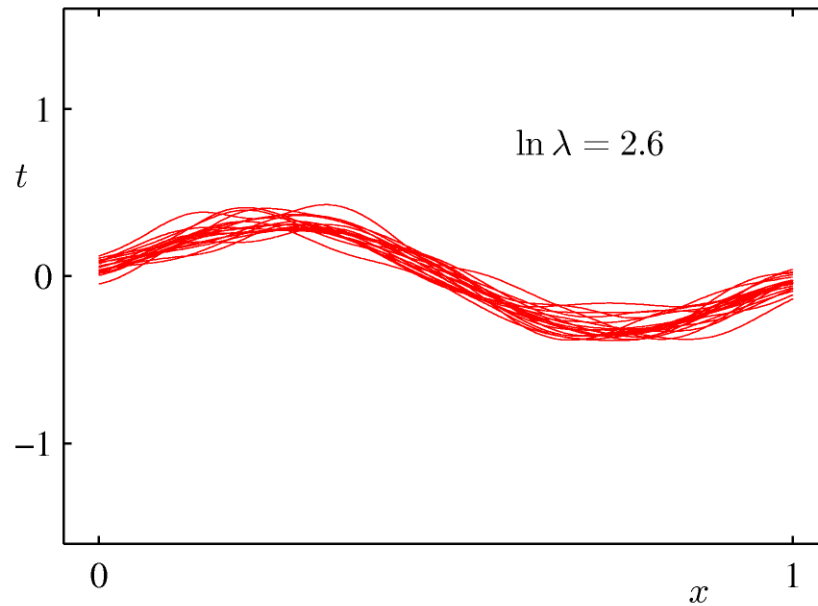
$$\textit{variance} = \int \mathbb{E}_{\mathcal{D}}[\{y(x; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(x; \mathcal{D})]\}^2] p(x) dx$$

$$\textit{noise} = \iint \{h(x) - t\}^2 p(x, t) dx dt$$



The Bias-Variance Decomposition₍₅₎

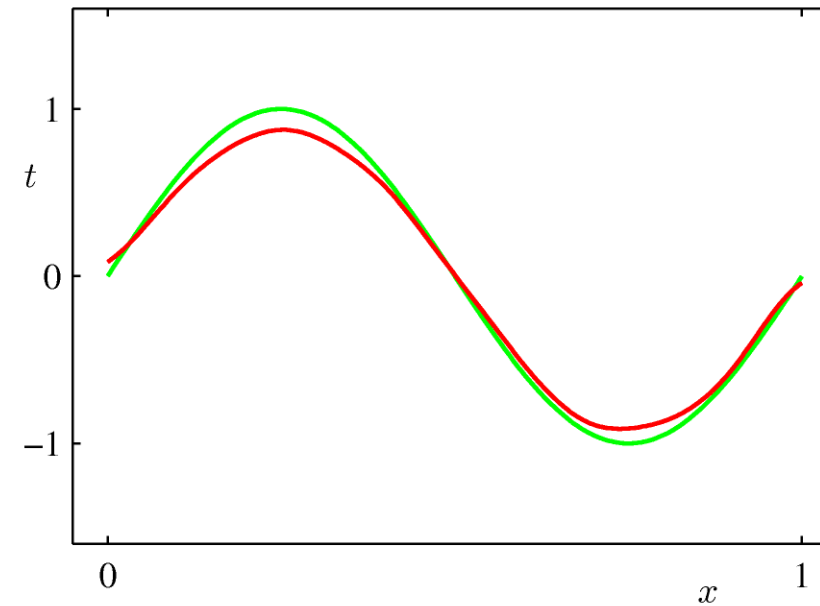
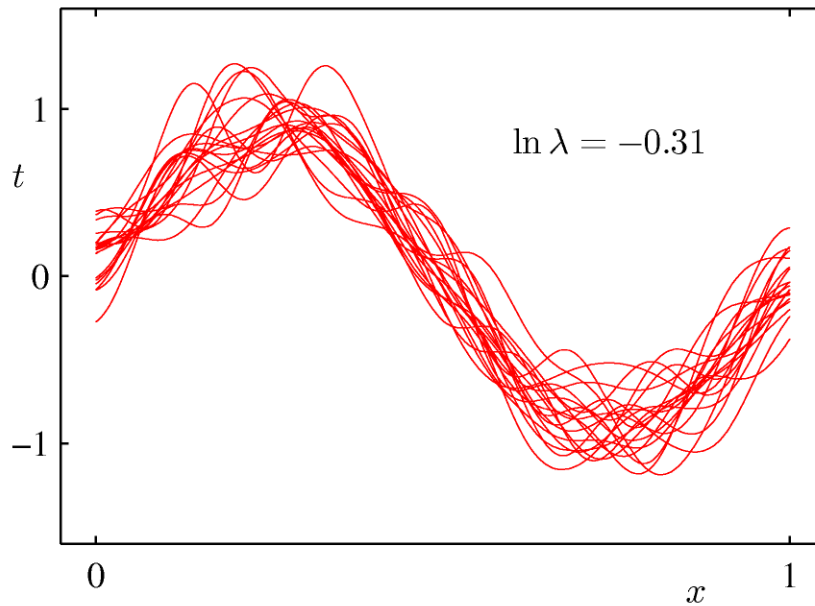
- Example: 25 data sets from the sinusoidal, varying the degree of regularization,





The Bias-Variance Decomposition₍₆₎

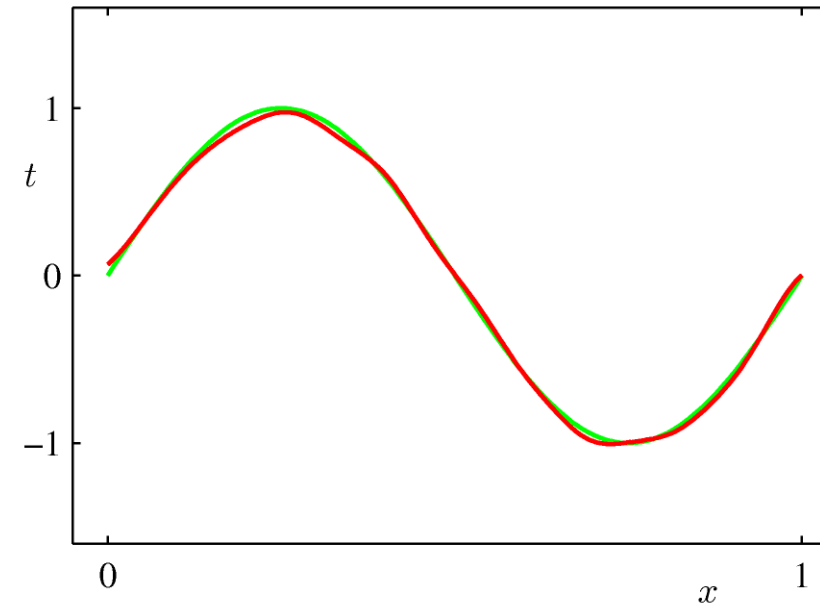
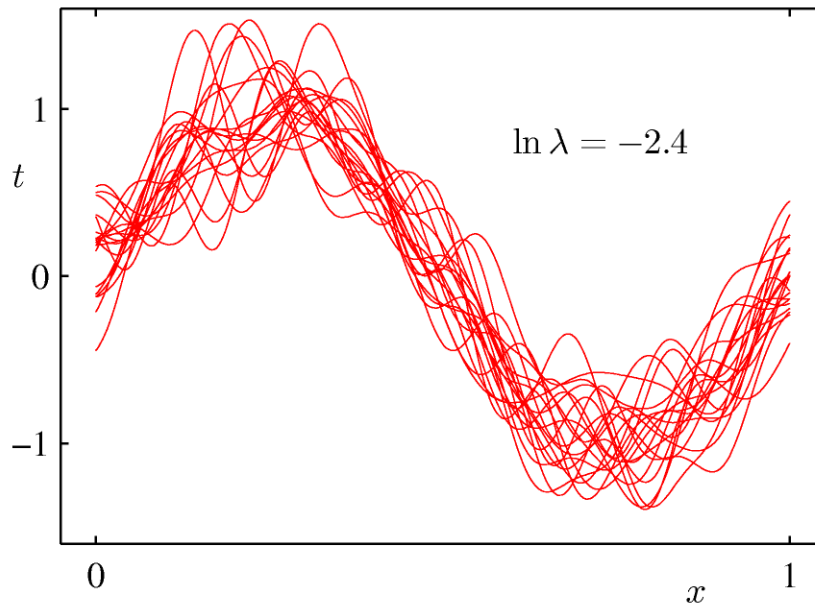
- Example: 25 data sets from the sinusoidal, varying the degree of regularization,





The Bias-Variance Decomposition₍₇₎

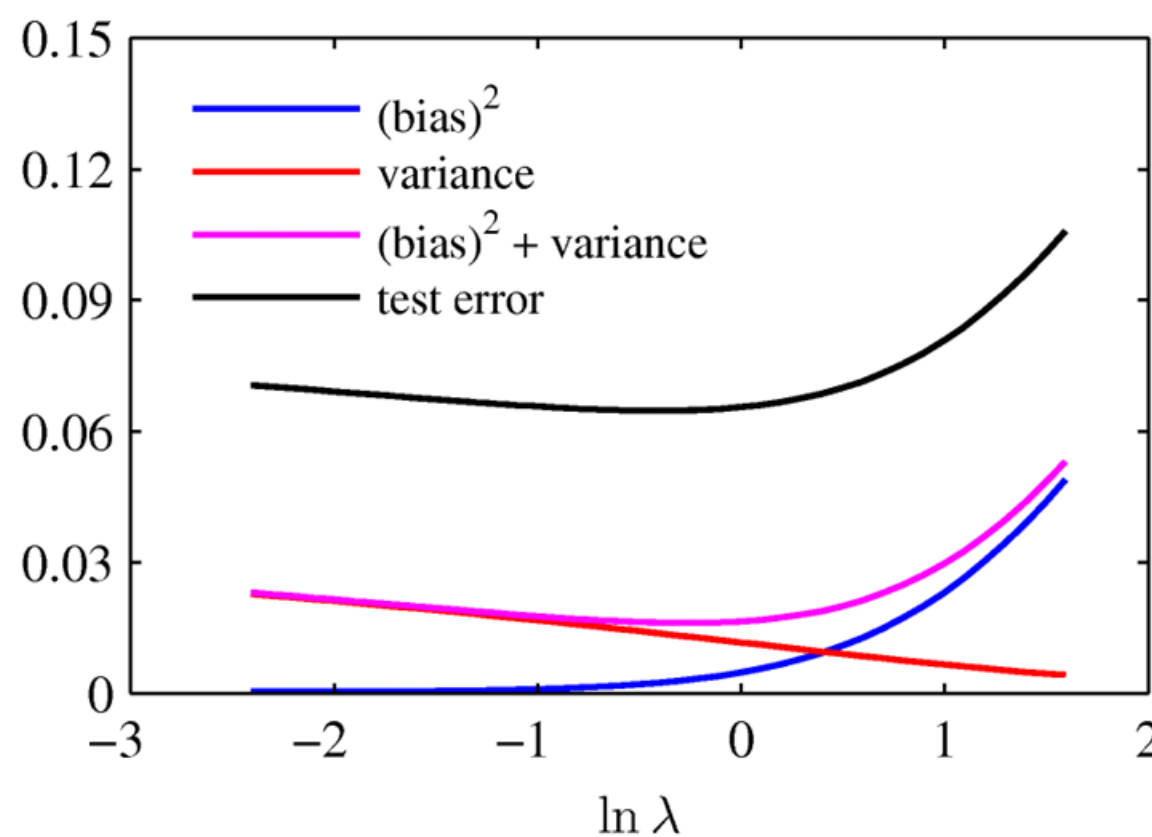
- Example: 25 data sets from the sinusoidal, varying the degree of regularization,





The Bias-Variance Trade-off

- From these plots, we note that an over-regularized model (large λ) will have a high bias, while an under-regularized model (small λ) will have a high variance.





Bayesian Linear Regression₍₁₎

- Define a conjugate prior over w

$$p(w) = \mathcal{N}(w|m_0, S_0)$$

Combining this with the likelihood function and using results for marginal and conditional Gaussian distributions, gives the posterior

where

$$p(w|t) = \mathcal{N}(w|m_N, S_N)$$

$$m_N = S_N(S_0^{-1}m_0 + \beta\Phi^T t)$$

$$S_N^{-1} = S_0^{-1} + \beta\Phi^T\Phi$$



Bayesian Linear Regression₍₂₎

- A common choice for the prior is

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I})$$

for which

$$m_N = \beta \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t}$$

$$\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}$$

Next we consider an example ...



Predictive Distribution₍₁₎

- Predict t for new values of x by integrating over w :

$$p(t|\mathbf{t}, \alpha, \beta) = \int p(t|\mathbf{w}, \beta) p(\mathbf{w}|\mathbf{t}, \alpha, \beta) d\mathbf{w}$$

where

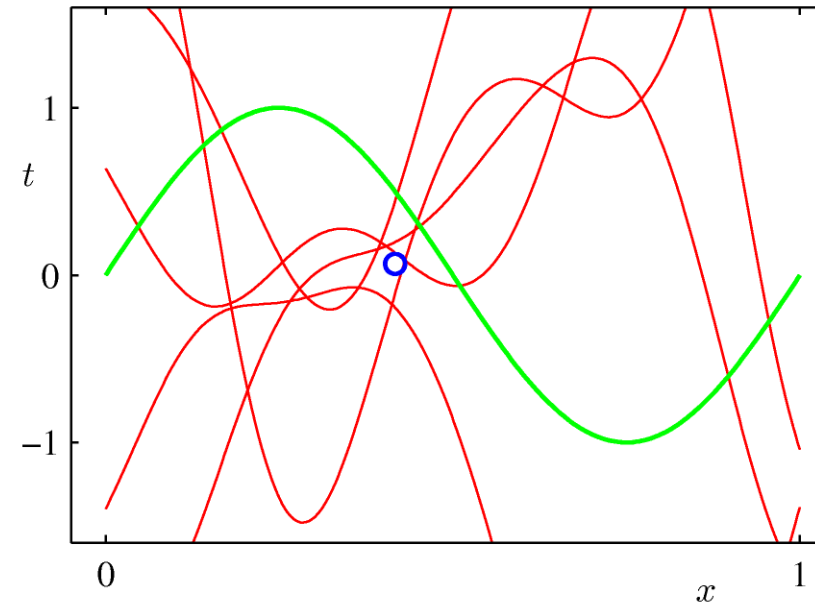
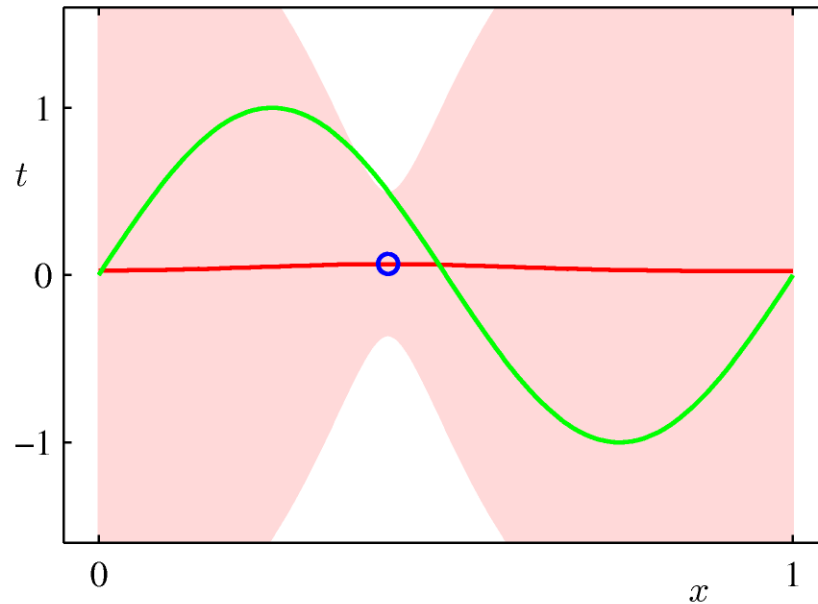
$$= \mathcal{N}(t|\mathbf{m}_N^T \phi(x), \sigma_N^2(x))$$

$$\sigma_N^2(x) = \frac{1}{\beta} + \phi(x)^T S_N \phi(x)$$



Predictive Distribution₍₂₎

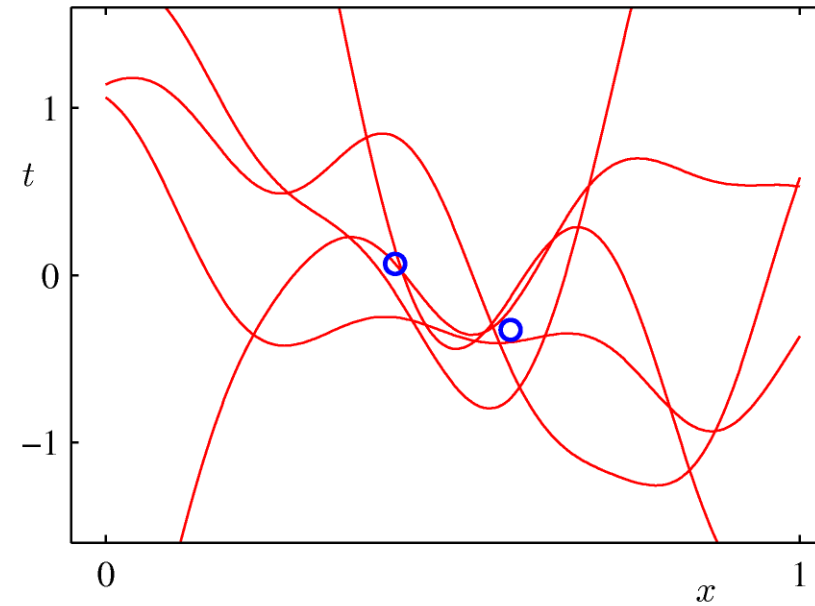
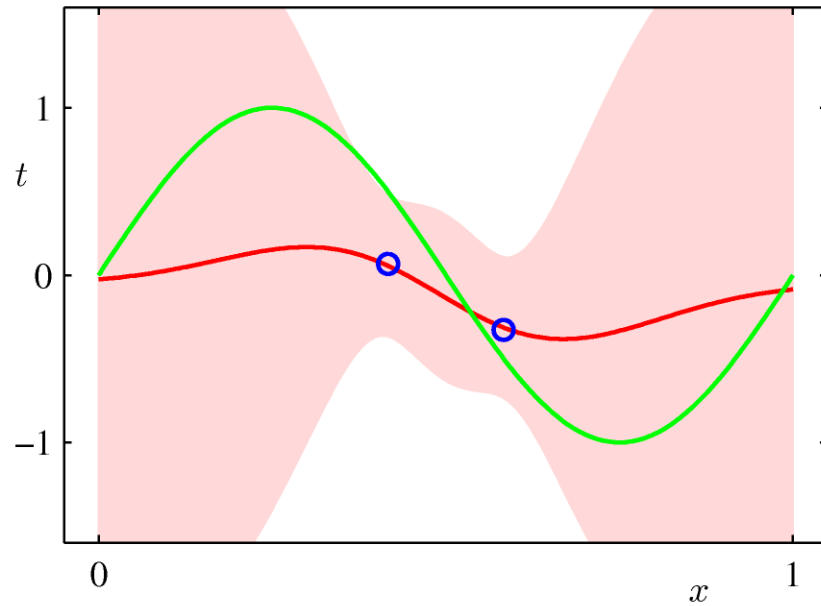
- Example: Sinusoidal data, 9 Gaussian basis functions, 1 data point





Predictive Distribution₍₃₎

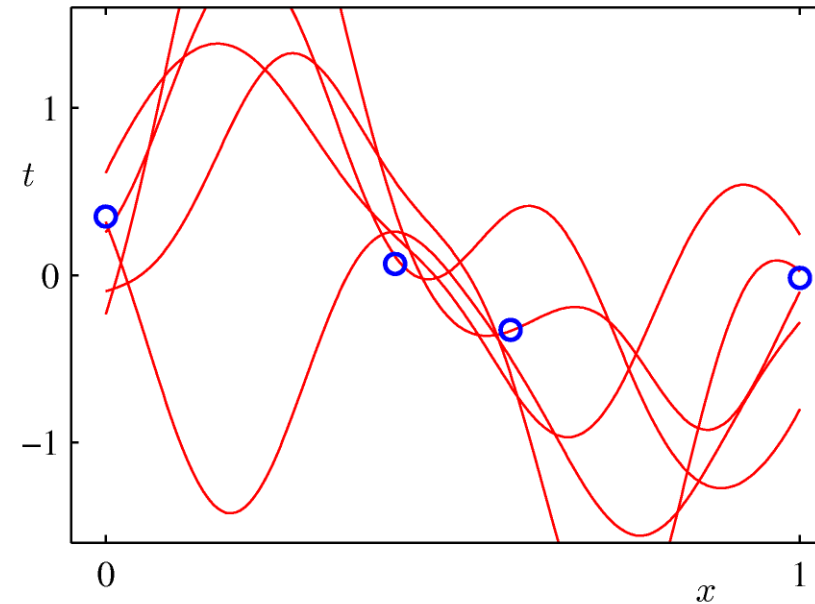
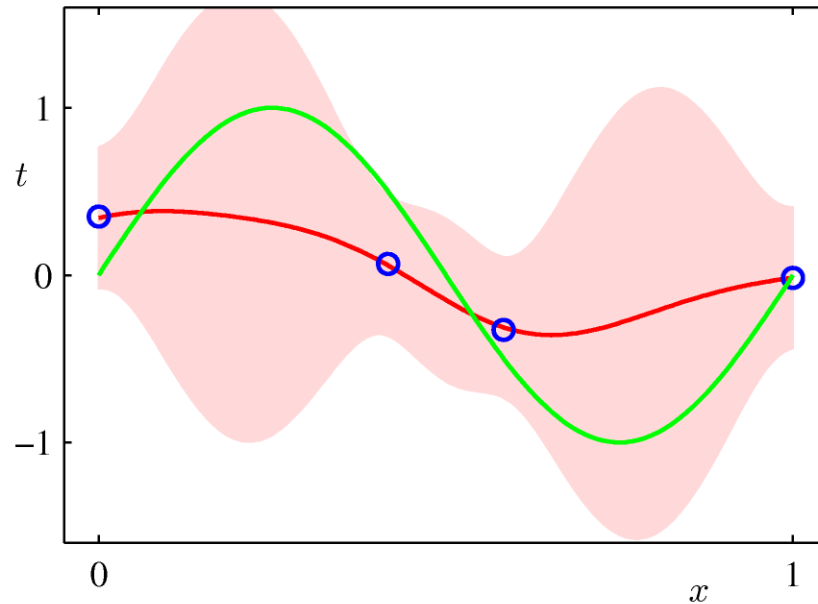
- Example: Sinusoidal data, 9 Gaussian basis functions, 2 data point





Predictive Distribution₍₄₎

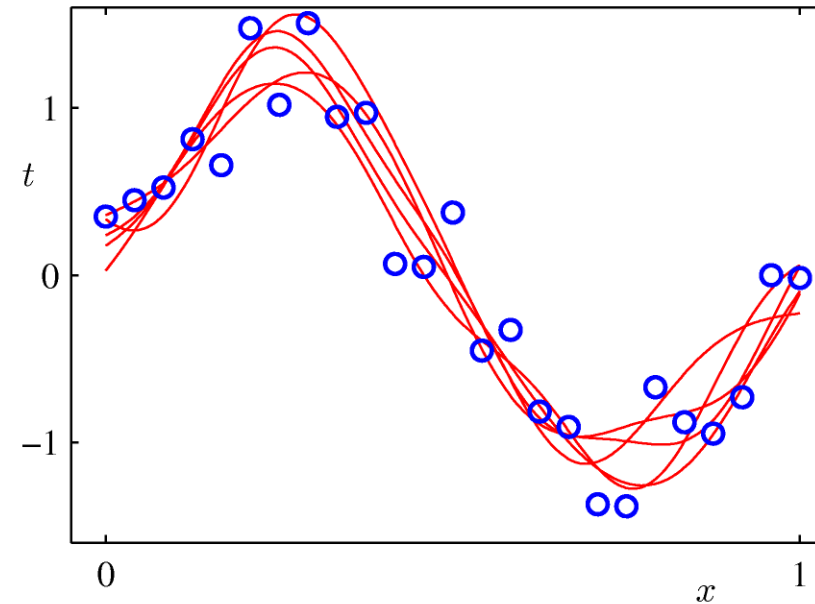
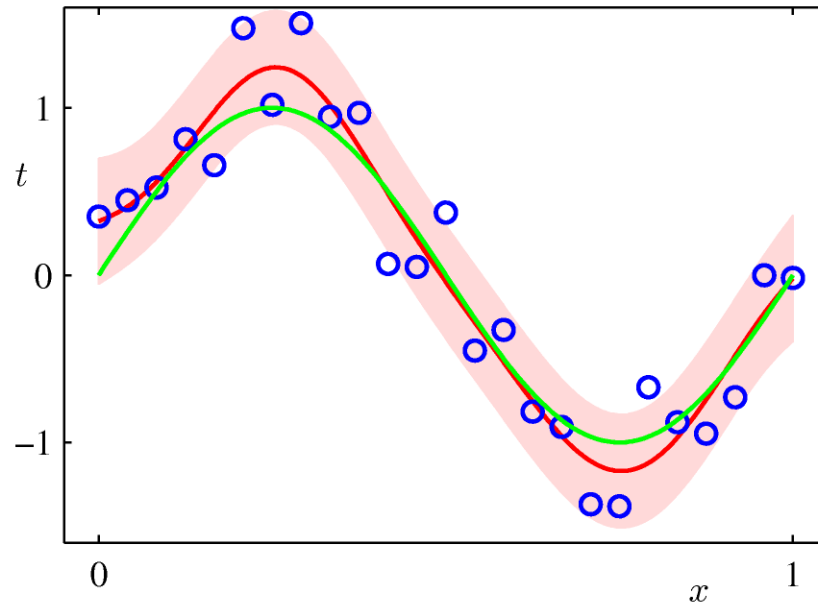
- Example: Sinusoidal data, 9 Gaussian basis functions, 4 data point





Predictive Distribution₍₄₎

- Example: Sinusoidal data, 9 Gaussian basis functions, 25 data point



Any questions?

