

Linear Models for Regression

Key concepts:

- Sum of squared error for regression loss
- Gradient descent for loss minimization
- Closed form solution for minimizing SSE
- Maximum likelihood estimation and SSE's probabilistic interpretation
- Overfitting and regularization

Example Regression Problems

- Predict housing price based on
 - House size, lot size, Location, # of rooms ...
- Predict stock price based on
 - Price history of the past month ...
- Predict the abundance of a species based on
 - Environmental conditions, shrubs? trees? ...

A Basic Set Up

Given: a training set, containing a total of N training examples (\mathbf{x}_i, y_i) , for $i = 1, \dots, N$

Goal: learn a function \hat{y} to minimize some loss function on the training data

To begin, we consider a linear hypothesis space (thus the name linear regression):

$$\text{Let } \mathbf{x} = [1, x_1, x_2, \dots, x_d]^T$$

$$\hat{y}(\mathbf{x}) = w_0 + w_1 x_1 + \dots + w_d x_d = \mathbf{w}^T \mathbf{x}$$

$$\text{where } \mathbf{w} = [w_0, w_1, \dots, w_d]^T$$

A Popular Loss Function: Sum-Squared Error

- Given a set of training examples

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2) \dots, (\mathbf{x}_N, y_N)$$

- Hypothesis space (representation of the function):

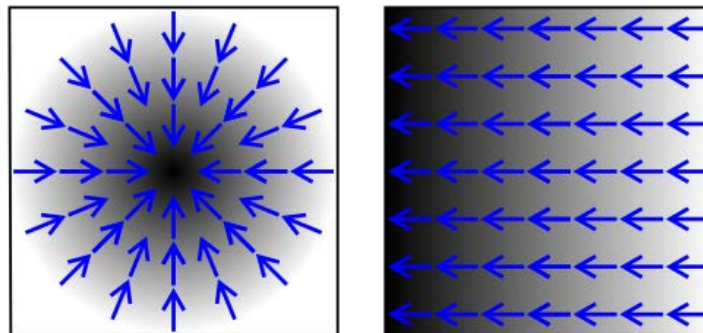
$$\hat{y}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$$

- Loss function (objective):

$$L(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i - y_i)^2$$

Optimization of the loss function

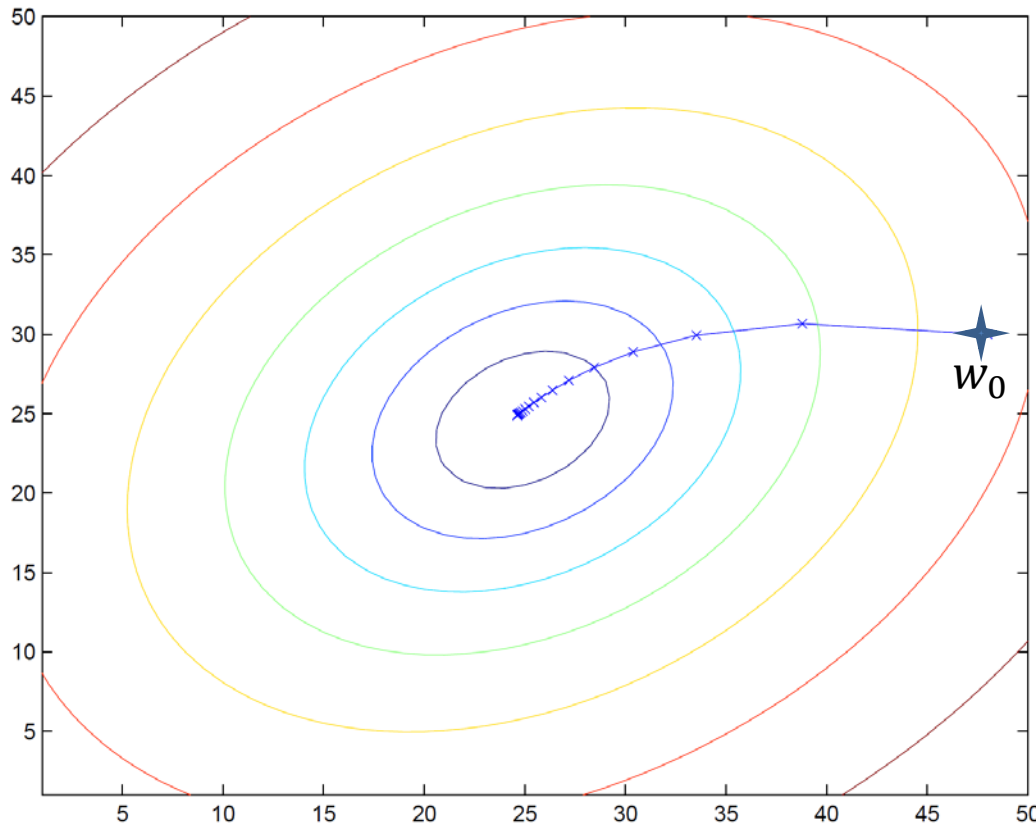
- Many ways to optimize this objective function
- A simple approach is **gradient descent**
 - Start with some random guess of the parameter
 - Iteratively improve the parameter by following the **steepest descent direction**
- **Gradient**: multivariate generalization of derivative, points in the direction of greatest rate of increase of the function



From Wikipedia

Values of the function in black and white, black representing higher values
Gradient represented by blue arrows.

Gradient descent to minimize $f(\mathbf{w})$



1. Start from some initial guess \mathbf{w}^0
2. Find the direction of steepest descent – opposite of the gradient direction $\nabla f(\mathbf{w})$
3. Take a step toward that direction (step size controlled by λ)
$$\mathbf{w}^{t+1} = \mathbf{w}^t - \lambda \nabla f(\mathbf{w}^t)$$
4. Repeat until no local improvement is possible
($|\nabla f(\mathbf{w}^t)| \leq \epsilon$)

Figure shows the contours of a quadratic function. Gradient descent starts from w_0 , and gradually moves toward the optimal solution.

From derivative to gradient

Derivative with examples:

$$f(x) = x^k; \quad \frac{\partial f}{\partial x} = kx^{k-1}$$
$$f(x) = e^x + 2x^2; \quad \frac{\partial f}{\partial x} = e^x + 4x$$

Generalizing to multivariate function:

For $f(\mathbf{x})$ where $\mathbf{x} = [x_1, \dots, x_d]^T$

$$\text{Gradient: } \nabla f(\mathbf{x}) = \left[\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_d} \right]$$

*Gradient is often represented as a row vector, but sometimes also as column vector for convenience of representation. Watch for context to ensure proper interpretation.

Read Chapter 5 Section 1-3 of Mathematics for Machine learning for the necessary background

Gradient Descent for SSE

$$L(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i - y_i)^2$$

We need to compute the gradient of $L(\mathbf{w})$:

$$\frac{\partial L}{\partial w_k} = \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i - y_i) x_{ik}$$

$$\nabla L(\mathbf{w}) = \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i - y_i) \mathbf{x}_i = \sum_{i=1}^N (\hat{y}_i - y_i) \mathbf{x}_i$$

Batch Gradient Descent for SSE

$$\mathbf{w} = \mathbf{w}^0$$

Repeat{

$$\nabla L(\mathbf{w}) = \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i - y_i) \mathbf{x}_i$$

$$\mathbf{w} \leftarrow \mathbf{w} - \lambda \nabla L(\mathbf{w})$$

} until $|\nabla L(\mathbf{w})| \leq \epsilon$

- The update rule goes through all training examples before taking each update step – thus called **batch gradient descent**
 - Examples with larger error $(\mathbf{w}^T \mathbf{x}_i - y_i)$ contribute more to the update
- **λ : learning rate or step size**
 - Large λ makes faster progress, but could lead to oscillation
 - Small λ avoids overstepping, but converges slower
- ϵ : convergence criterion
- The SSE Objective is convex – converge to global optimal

Stochastic gradient descent

$$\mathbf{w} = \mathbf{w}^0$$

Repeat until convergence {

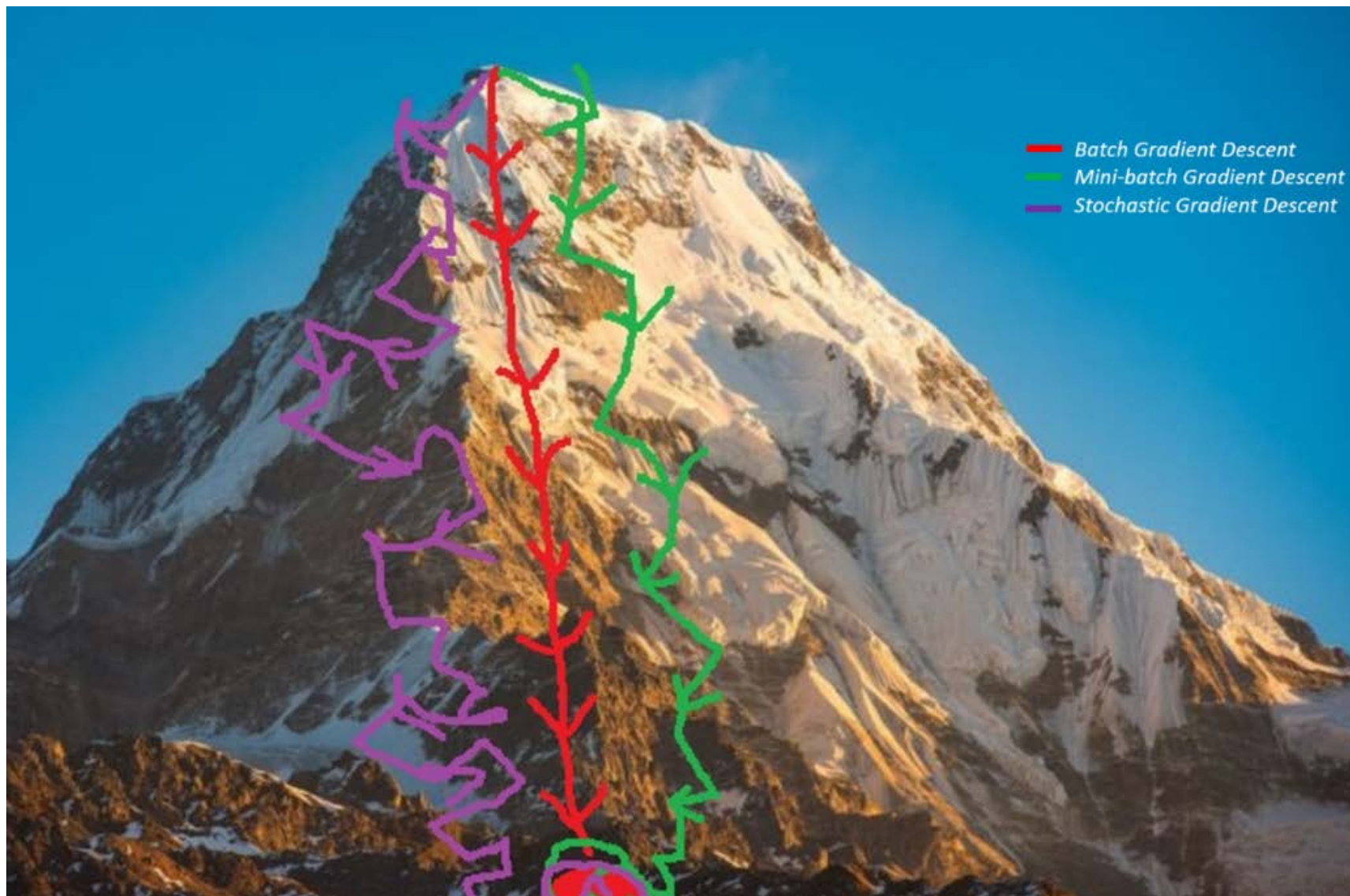
$$\text{for } i = 1 \text{ to } N \quad \mathbf{w} \leftarrow \mathbf{w} - \lambda(\mathbf{w}^T \mathbf{x}_i - y_i) \mathbf{x}_i$$

}

- It repeatedly runs through the training examples
- It updates \mathbf{w} immediately when seeing an example, based on the gradient w.r.t. that training instance

Batch, Minibatch, Stochastic gradient descent

- Batch: go through all training examples for computing the gradient
 - not efficient sometimes infeasible for large data
- Stochastic gradient descent: compute gradient contribution for each training example and update immediately
 - Training can be jumpy
- Minibatch: go through a subset of the training example before updating
 - Can be a good tradeoff



Alternative way to optimizing $L(\mathbf{w})$

- By setting the gradient to zero, we can directly find the extreme (minimum) point of $L(\mathbf{w})$
- Define:

$$\mathbf{X} = \begin{bmatrix} x_{10} & x_{11} & \cdots & x_{1d} \\ x_{20} & x_{21} & \cdots & x_{2d} \\ \cdots & \cdots & \cdots & \cdots \\ x_{N0} & x_{N1} & \cdots & x_{Nd} \end{bmatrix}, \quad \mathbf{Y} = \begin{bmatrix} y_1 \\ \cdots \\ y_N \end{bmatrix}$$

We can show that

$$\nabla E(\mathbf{w}) = \mathbf{X}^T(\mathbf{X}\mathbf{w} - \mathbf{Y})$$

Setting it to zero:

$$\Rightarrow \mathbf{X}^T\mathbf{X}\mathbf{w} - \mathbf{X}^T\mathbf{Y} = 0$$

$$\Rightarrow \mathbf{X}^T\mathbf{X}\mathbf{w} = \mathbf{X}^T\mathbf{Y} \quad \Rightarrow \quad \mathbf{w} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y}$$

Normal equation

Probabilistic Interpretation of SSE

- One might ask, why is SSE a good objective?
- There is a natural probabilistic explanation for SSE with some simple probabilistic assumptions
- Before that, we will provide a brief of maximum likelihood estimation

Background: Maximum Likelihood Estimation

Parameter Estimation

- Given observation of some random variable(s), assuming the variable(s) follow some fixed distribution, how to estimate the parameters?
- Example:
 - coin tosses (Bernoulli distribution, parameter: p – *probability of head*)
 - Dice or grades (Discrete distribution, parameters: p_i for $i = 1, \dots, k - 1$ with k categories)
 - Height and weight of a person (continuous variable, parameters depends on the distribution, for example for Gaussian Distribution, the parameters are μ , the mean and Σ , the covariance)

Maximum Likelihood Principle

- We will use a general term M to denote the model (parameters) we try to estimate
- We will use a general term D to denote the data that we observe
 - D is a set of examples, let D_m denote the m -th example in D
- Assuming a fixed model M , what is the probability of observing D ?
 - This can be written as $P(D; M)$
- The maximum likelihood principle seeks to find M that maximizes this probability
- This is called the likelihood function
 - $P(D; M) = \prod_{m=1}^n P(D_m; M)$ --- here we are making the IID assumption, that is examples in D are **independently and identically distributed**, thus we can use the product
- It is often more convenient to work with the log of $P(D; M)$
 - $l(M) = \log P(D; M) = \sum_{m=1}^n \log P(D_m; M)$

Example: Coin Toss



- You are given a coin, and want to decide p - the probability of head of this coin
- You toss it for 500 times, and observe heads 242 times
- What is your estimate of p ?

$$p = \frac{242}{500} = 0.484$$

- But why? – this is doing maximum likelihood estimation
- In other words, $p = 0.484$ leads to the highest probability of observing 242 heads out of 500 tosses among all possible values of p
- Now let's go over the math to convince ourselves

Example Cont.

- We have $D = \{x_1, \dots, x_N\}$, where x_i is the outcome of the m -th coin toss,
 - 1 means head, and 0 means tail
- Lets write down the likelihood function: $P(D; p) = \prod_{i=1}^N P(x_i; p)$
- For binary variable, we have a nice compact form to write down its probability mass function
$$P(x_i; p) = p^{x_i}(1 - p)^{1-x_i}$$
- So we have

$$P(D; p) = \prod_{i=1}^N p^{x_i}(1 - p)^{1-x_i}$$

- Take the log:

$$\begin{aligned} l(p) &= \log P(D; p) = \sum_{i=1}^N \log p^{x_i}(1 - p)^{1-x_i} \\ &= \underbrace{\sum_{i=1}^N x_i}_{n_1} \log p + \underbrace{\sum_{i=1}^N (1 - x_i)}_{n_0} \log(1 - p) = n_1 \log p + n_0 \log(1 - p) \end{aligned}$$

Maximizing the likelihood

- Take derivative of $l(p) = n_1 \log p + n_0 \log(1 - p)$

$$\frac{dl(p)}{dp} = \frac{n_1}{p} - \frac{n_0}{1 - p}$$

- Setting it to zero:

$$\frac{n_1}{p} = \frac{n_0}{1 - p}$$

- Solving this leads to:

$$p = \frac{n_1}{n}$$

End of quick
intro on MLE

Back to Linear Regression

- **Assumption:** target variable y and the input \mathbf{x} are related as follows:

$$y = \mathbf{w}^T \mathbf{x} + \epsilon$$

where $\epsilon \sim N(0, \sigma^2)$, i.e., it is Gaussian noise

- Equivalently we have: $y|\mathbf{x} \sim N(\mathbf{w}^T \mathbf{x}, \sigma^2)$
- Given a set of observed \mathbf{x} and y pairs: $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2) \dots, (\mathbf{x}_N, y_N)$, **independently and identically distributed (IID)**, use maximum likelihood estimation to estimate \mathbf{w}
- **Likelihood function:** given the inputs $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$, what is the probability of observing target outputs y_1, y_2, \dots, y_N ?

$$P(y_1, \dots, y_N | \mathbf{x}_1, \dots, \mathbf{x}_N; \mathbf{w}) = \prod_{i=1}^N P(y_i | \mathbf{x}_i; \mathbf{w})$$

Maximum likelihood estimation of \mathbf{w}

We will now work with log of the likelihood:

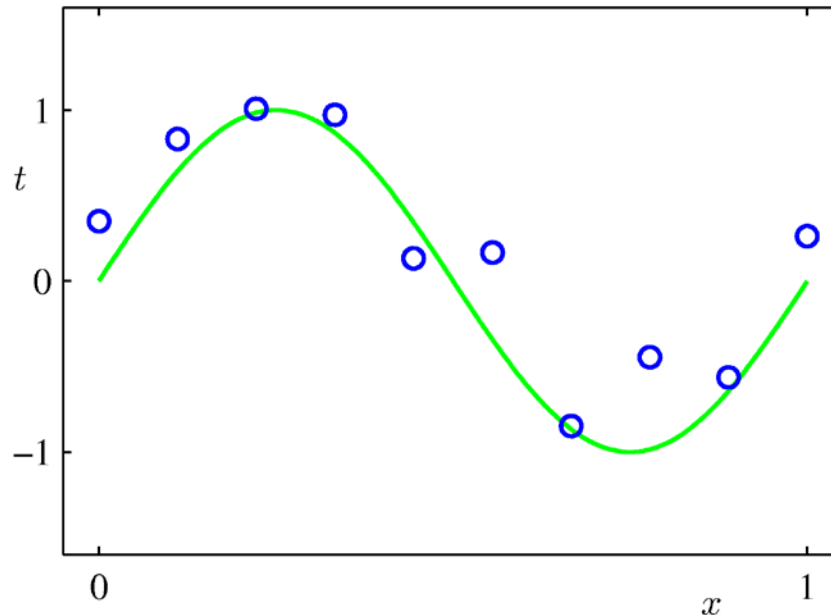
$$\begin{aligned} l(\mathbf{w}) &= \log \prod_{i=1}^N P(y_i | \mathbf{x}_i; \mathbf{w}) \\ &= \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi}\sigma} \exp \frac{-(y_i - \mathbf{w}^T \mathbf{x}_i)^2}{2\sigma^2} \\ &= -N \log \sqrt{2\pi}\sigma + \sum_{i=1}^N \frac{-(y_i - \mathbf{w}^T \mathbf{x}_i)^2}{2\sigma^2} \\ \arg \max_{\mathbf{w}} l(\mathbf{w}) &= \boxed{\arg \min_{\mathbf{w}} \frac{1}{2} \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2} \end{aligned}$$

Maximum Likelihood Estimation under Gaussian noise leads to
Sum of Squared Error objective

What we have seen so far

- We introduced linear regression
 - Goal: learn a linear function $\hat{y}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$
 - Approach:
 - Choose a loss function: SSE
 - Optimize the loss function (Gradient descent, or analytically solving the optimization problem)
- What if the relationship between \mathbf{x} and y is not linear? For example, quadratic or cubic?
 - Linear regression can still be used
 - All we need is to introduce nonlinear features

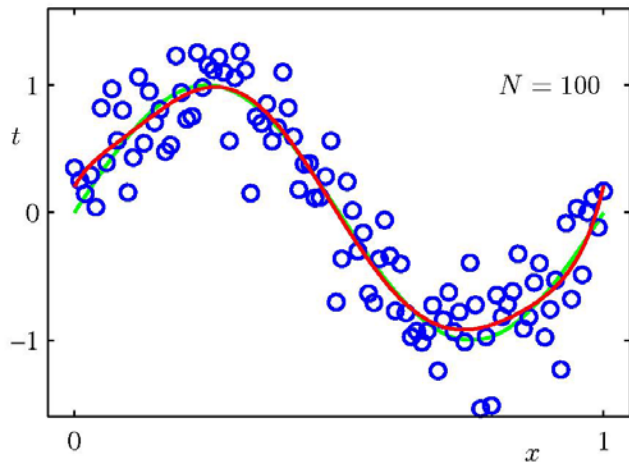
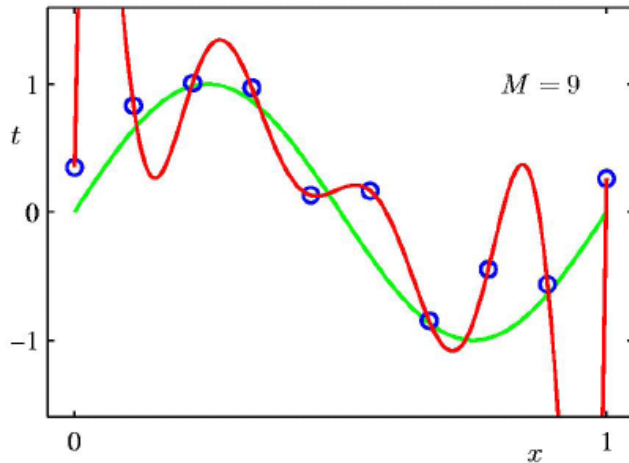
- Revisit example: Polynomial Curve Fitting



$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_jx^j$$

- In this example, there is only one feature x . We learn a function of M -order polynomial
- Alternatively, we could also view this as learning a linear model considering $(1, x, x^2, \dots, x^M)$ as the features.
- Note that this new feature space is derived from the original input x
- We refer to such derived features as **the basis functions**

Over-fitting issue



- What can we do to curb over-fitting
 - Use less complex model
 - Use more training examples
 - **Regularization**

In linear regression, overfitting can often be characterized by large weights

	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

Regularized Linear Regression

- Consider the following loss function:

$$L_D(\mathbf{w}) + \lambda L_w(\mathbf{w})$$

Data term + Regularization term (penalize complex models)

$$\sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + \lambda \sum_{j=1}^M |w_j|^q$$

Encourage small weight values

	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

L2 Regularized Linear Regression

- With the SSE loss and a **quadratic regularizer**, we get

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

- which is minimized by

$$\mathbf{w} = (\lambda \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

- λ : regularization coefficient, which controls the trade-off between model complexity and the fit to the data
 - Larger λ encourages simple model (driving more elements of \mathbf{w} to 0)
 - Small λ encourages better fit of the data (driving SSE to zero)

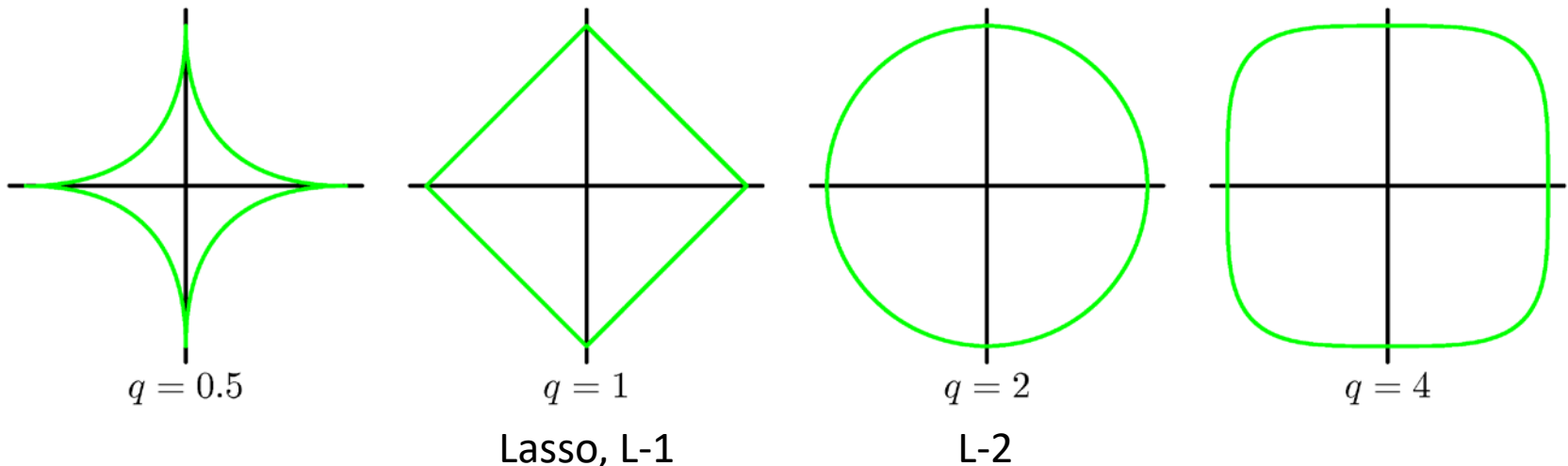
More Regularizations

$$\sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + \lambda \sum_{j=0}^M |w_j|^q$$

Equivalent to minimizing SSE subject to $\sum_{i=0}^M |w_i|^q \leq \epsilon$

A good explanation of this equivalence is provided here:

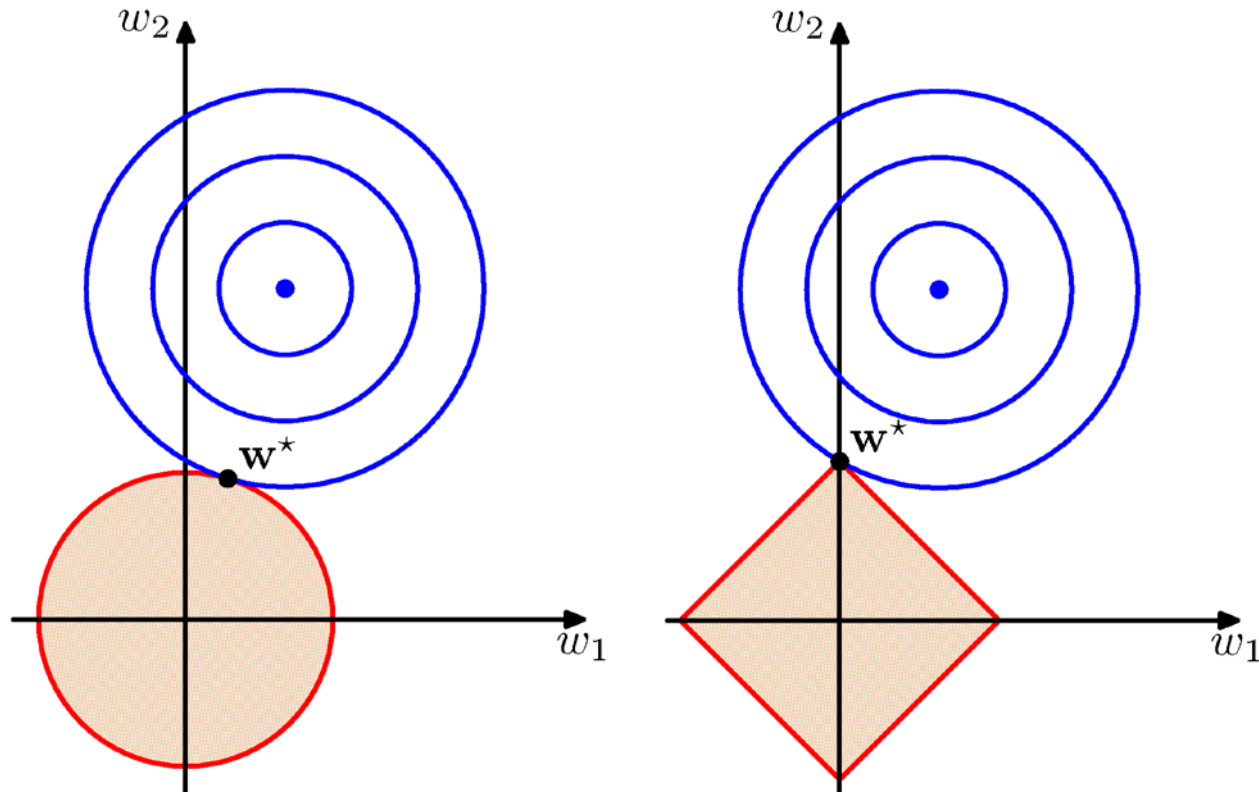
<http://math.stackexchange.com/questions/335306/why-are-additional-constraint-and-penalty-term-equivalent-in-ridge-regression>



Shape is determined by q , size determined by λ

Regularized Linear Regression

- Lasso ($q = 1$) tends to generate sparser solutions (majority of the weights shrink to zero) than a quadratic regularizer ($q = 2$, often called ridge regression).



Commonly used regularizers

- L-2 regularization
$$\sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + \lambda \sum_{j=0}^M w_j^2$$

Poly-time close-form solution
Curbs overfitting but does not produce sparse solution

- L-1 regularization
$$\sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + \sum_{j=0}^M |w_j|$$

Poly-time approximation algorithm
Sparse solution – potentially many zeros in \mathbf{w}

- L-0 regularization
$$\sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + \sum_{j=0}^M I(w_j \neq 0)$$

Seek to identify optimal feature subset
NP-complete problem!

Summary

- Linear regression with SSE objective
 - Gradient descent algorithm for optimization
 - Closed-form solution by solving the normal equation
 - Maximum likelihood estimation assuming IID Gaussian noise with zero mean
- Maximum likelihood estimation: a basic probabilistic principle for setting up learning objectives
- Using basis functions allows us to learn nonlinear functions
- Controlling over-fitting by Regularization
 - The regularization term is equivalent to adding a constraint to the parameters
 - In effect, Quadratic regularization for linear regression adds a diagonal matrix to $X^T X \rightarrow (\lambda I + X^T X)$ – it also serve the purpose of conditioning $X^T X$ when it is not full rank (not invertible)
 - LASSO regularization encourages sparse solution