

Regularization

JrPhy

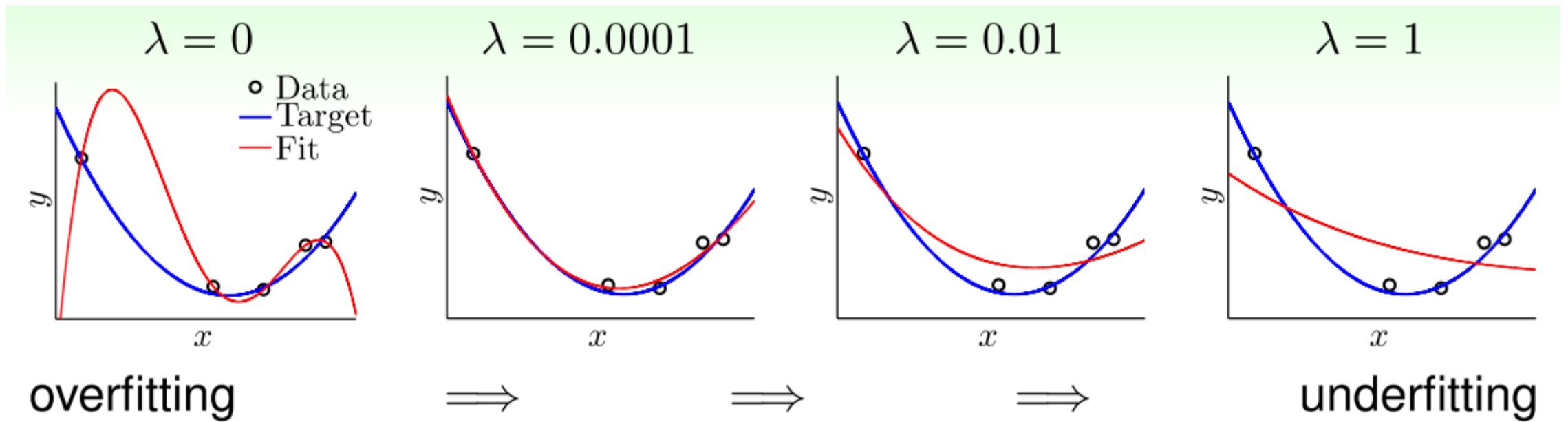
Introduction

- So far we can accept some mistakes, but we still want to the mistakes are the least, it's called regularization. Now we've find a hyperplane and add a constrain $\min_w E_{in}(w)$ such that $\sum_i w_i^2 \leq C$
- This problem can be solve by Lagrange undetermined multiplier, here we use square error, let $z_i = f(x_i)$ is our transformation, and the error is
- $E(w_i) = (z_i w^T - y_i)^2$, $E = \sum_i E(w_i) = (Zw^T - Y)(Zw^T - Y)^T$,

Optimize

- What we want to solve is $\min_w \left((Zw^T - Y)(Zw^T - Y)^T + \frac{\lambda}{N} ww^T \right) \quad \lambda > 0$
$$\frac{\partial}{\partial w} \left((Zw^T - Y)(Zw^T - Y)^T + \frac{\lambda}{N} ww^T \right) = 0$$
$$\frac{2}{N} (ZZ^T w^T - Z^T Y) + \frac{2\lambda}{N} w^T = 0 \rightarrow w = (ZZ^T + \lambda I)^{-1} ZY$$
- As the λ is bigger, then higher order term increase fast, so its coefficient is usually small. So that λ can suppress higher order term as it increases a small number

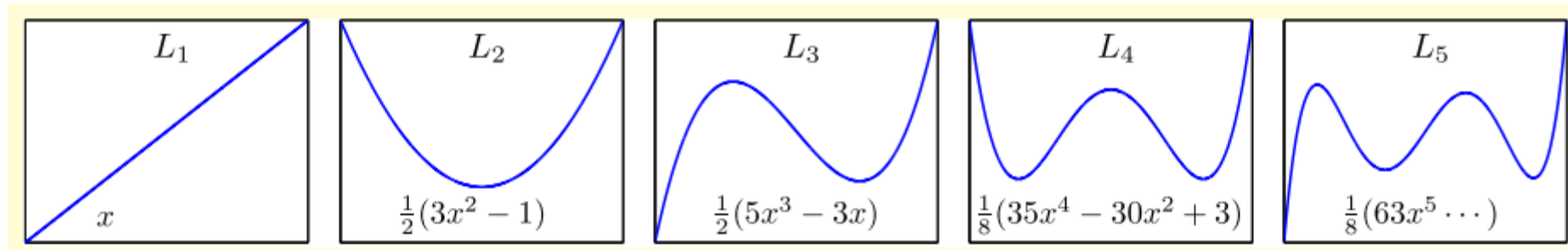
Choosing λ



A tip of polynomial transform

- We've mentioned previous before, if you use polynomial transform, it is recommended to use Legendre transform, it's also a polynomial but with better numerical property.

- $$-1 \leq L(x_i) \leq 1, \quad \int_{-1}^1 L(x_i)L(x_j)dx = \delta_{ij}$$



L1 and L2 regularization

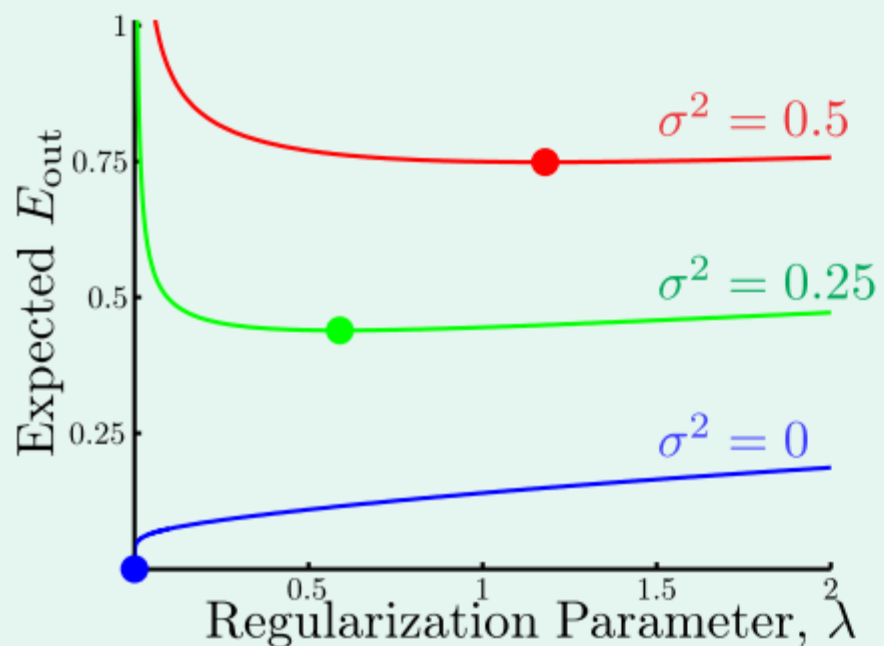
- So far we use the square error and we want to minimize it, it's called the L2 regularization in machine learning. There is another regularization, L1 regularization, its form is $|\mathbf{w}|$, both regularizations can avoid the noise.
- L2 regularization " $\mathbf{w}\mathbf{w}^T$ ": easy to optimize, differential everywhere
- L1 regularization " $|\mathbf{w}|$ ": hard to optimize, not differential at $|\mathbf{w}|$.
- The optimal solution with L1 regularization is near $|\mathbf{w}| = 0$, so the most coefficients are 0, it's called "sparse solution", so it just needs calculate a few terms.

Use your own regularization

- It's ok that choose your own regularization, there are some criterions of the choices:
- 1. target dependent
- 2. what the user want
- 3. easy to optimize
- If your regularization is bad, λ can protect your choice.

Optimal λ

stochastic noise



deterministic noise

