

# Supervised Regression

## Design matrix

Let  $x_1, \dots, x_p$  be our  $p$  feature vectors, then we call the matrix  $X = \begin{pmatrix} 1 & x_1 & \dots & x_p \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & \dots & x_p \end{pmatrix} \in \mathbb{R}^{n \times (p+1)}$  the **design matrix**.  
Furthermore, the define  $x = (1, x_1, \dots, x_p) \in \mathbb{R}^{p+1}$  and  $\theta = (\theta_0, \theta_1, \dots, \theta_p) \in \mathbb{R}^{p+1}$   
one row as col vector (X has n of those rows), smth  $\tilde{X} = (1, X)$  for clarity to distinguish intercept / feature ver.

## Setup

We predict  $y \in \mathbb{R}$  as linear combinations of features  
 $\hat{y} = f(x) = \theta^T x = \theta_0 + \theta_1 x_1 + \dots + \theta_p x_p$   
This results in the hypothesis space  $\mathcal{H} = \{f(x) = \theta^T x \mid \theta \in \mathbb{R}^{p+1}\}$   
A typical choice for the Loss-function is L2-Loss.  
This results in  $R_{emp}(\theta) = \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2 = SSE = \sum_{i=1}^n \frac{1}{2} (y^{(i)} - \theta^T x^{(i)})^2 = \frac{1}{2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2$

## Optimization

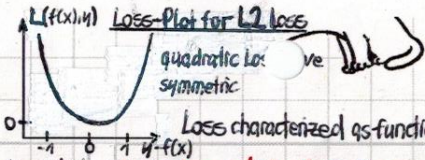
We want to find  $\hat{\theta} = \arg \min_{\theta \in \mathbb{R}^{p+1}} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2 = \arg \min_{\theta \in \mathbb{R}^{p+1}} \|y - X\theta\|_2^2$   
This results in the ordinary-least-squares estimator:  
 $\hat{\theta} = (X^T X)^{-1} X^T y$  (MLE) **Gradient:**  $\nabla_{\theta} R_{emp}(\theta) = \frac{2}{n} (X^T X \theta - X^T y) = 2(X^T X \theta - X^T y)$

### ANALYTICAL OPTIMIZATION - PROOF

$R_{emp}(\theta) = \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2 = \|y - X\theta\|_2^2$ ;  $\theta \in \mathbb{R}^p$  with  $p := p+1$   
using L2 loss

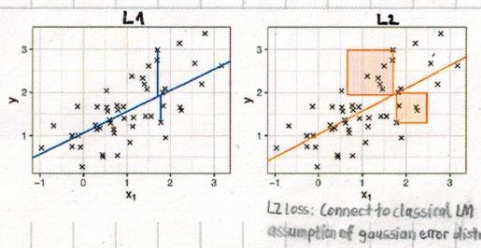
|  |   |
|--|---|
| $0 = \frac{\partial R_{emp}(\theta)}{\partial \theta}$ (sum notation)  | $0 = \frac{\partial R_{emp}(\theta)}{\partial \theta}$ (matrix notation)                      |
| $0 = \frac{\partial}{\partial \theta} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2$ sum & chain rule            | $0 = \frac{\partial \ y - X\theta\ _2^2}{\partial \theta}$                                    |
| $0 = \sum_{i=1}^n \frac{\partial}{\partial \theta} (y^{(i)} - \theta^T x^{(i)})^2$                             | $0 = \frac{\partial}{\partial \theta} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2$ chain rule |
| $0 = \sum_{i=1}^n 2(y^{(i)} - \theta^T x^{(i)}) \frac{\partial}{\partial \theta} (y^{(i)} - \theta^T x^{(i)})$ | $0 = \frac{\partial}{\partial \theta} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2$            |
| $0 = \sum_{i=1}^n 2(y^{(i)} - \theta^T x^{(i)}) (-x^{(i)})^T$  | $0 = \frac{\partial}{\partial \theta} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2$            |
| $0 = \sum_{i=1}^n -2(y^{(i)} - \theta^T x^{(i)}) x^{(i)}$  | $0 = \frac{\partial}{\partial \theta} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2$            |
| $0 = \sum_{i=1}^n -2(y^{(i)} - \theta^T x^{(i)}) x^{(i)}$  | $0 = \frac{\partial}{\partial \theta} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2$            |
| $0 = \sum_{i=1}^n -2(y^{(i)} - \theta^T x^{(i)}) x^{(i)}$  | $0 = \frac{\partial}{\partial \theta} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2$            |
| $0 = \sum_{i=1}^n -2(y^{(i)} - \theta^T x^{(i)}) x^{(i)}$  | $0 = \frac{\partial}{\partial \theta} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2$            |

encoding is basically One-Hot-Enc. with one reference, remaining = 0 (for categorical features)  
**Sidenote:** Linear Regression with intercept must use leave one out dummy encoding!  
One-Hot-Encoding (full k levels for k categories) leads to "over-parameterization" and thus an indistinguishable non-unique parameter vector, unstable possibly ridiculous  $\theta$  estimates  
but: **Regularized Regression** requires full encoding, otherwise intercept not penalized!



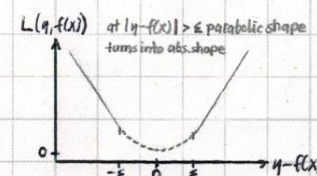
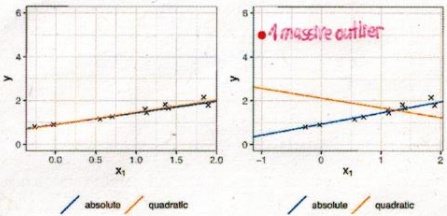
## L2-Loss / Squared Error

penalizes quadr. residuals  
 $L(y, f(x)) = \frac{1}{2} (y - f(x))^2$  or  $L(y, f(x)) = (y - f(x))^2$   
Error  $\perp$  Feature Space  $X^T (y - X\theta) = 0$   
Properties: Convex; differentiable everywhere  
Problem: Outliers have a huge effect on L2.



## L1-Loss / Absolute Error

penalizes absolute residuals  
 $L(y, f(x)) = |y - f(x)|$   
Properties: Convex; Robust against outliers  
Problem: Not differentiable, slower to optimize (only numerically)



## Huber Loss

Weighted piecewise combo of L1, L2 loss  
 $\epsilon$  marks where L2 transits to L1 loss ( $\epsilon > 0$ )  
 $L(y, f(x)) = \begin{cases} \frac{1}{2} (y - f(x))^2 & \text{if } |y - f(x)| \leq \epsilon \\ \epsilon |y - f(x)| - \frac{1}{2} \epsilon^2 & \text{else} \end{cases}$   
Huber loss combines advantages of both  
It's smooth & differentiable like L2 loss  
Doesn't punish large residuals harshly just as L1 loss  $\Rightarrow$  more robust

## Polynomial Regression

### General Linear Model

We predict  $y \in \mathbb{R}$  as linear combinations of basis functions  $\phi_j$   
 $\hat{y} = f(x) = \theta_0 + \theta_1 \phi_1(x_1) + \dots + \theta_p \phi_p(x_p)$   
If we set  $\phi_j = x_j$  we get the usual linear regression model.  
! This model is linear in parameters  $\theta$  but not in variables/features  $x$ .  
 $f(a x + b x^2; \theta) = a f(x; \theta) + b f(x^2; \theta)$  but  $f(x; a\theta + b\theta^2) = a f(x; \theta) + b f(x; \theta^2)$   
The design matrix now is defined as  $X = \begin{pmatrix} 1 & \phi_1(x_1) & \dots & \phi_p(x_p) \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \phi_1(x_n) & \dots & \phi_p(x_n) \end{pmatrix}$  can still optimize with normal equation

### Polynomial Regression

More flexible than simple Linear regression  
In the Polynomial Regression we choose the basis functions  $\phi: \mathbb{R} \rightarrow \mathbb{R}, x_i \mapsto \sum_{k=0}^d \beta_k x_i^k$   
This results in  $f(x) = \theta_0 + \theta_1 \phi(x) = \theta_0 + \sum_{k=1}^d \theta_k x^k$  where  $\theta_{k,k} = \theta_k \cdot \beta_k$   $\theta_k$  is just param. for full  $\phi(x)$   
here just one feature  $x$  idea: merge  $\theta_1$  and weight  $\beta_k$  (still linear in merged parameter)  
 $X = \begin{pmatrix} 1 & x^{(1)} & (x^{(1)})^2 & \dots & (x^{(1)})^d \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x^{(n)} & (x^{(n)})^2 & \dots & (x^{(n)})^d \end{pmatrix}, \theta \in \mathbb{R}^{d+1}$   
For two features we might model like this:  
 $f(x) = \theta_0 + \theta_1 x_1 + \sum_{k=2}^d \theta_{2,k} x_2^k$  polynomial for  $x_2$  up until grade  $d$   
 $x_1 = 1 - x_2 - x_3 - \dots - x_p = 0$ , only "hot"  $x_1 = 1 \Rightarrow (\beta_0, \beta_1, \dots, \beta_p)$  e.g. using  $x_1 = 1 - x_2$  linear term  $x_1$

### Overfitting risk with complex polynomials / basis functions

How high can / should we go with degree of modelling polynomials? (Hyperparameter!)  
Problem with high degree: Our model will fit the training data excellent, but generally will not perform well on new (test) data  
Data contains random noise that is not part of true data-generating process  
 $\Rightarrow$  model with overly high capacity learns all those spurious patterns and those will not reappear in new data, leading to poor generalization of our model  
Also higher degrees can lead to instability (oscillation esp. at bounds, "Runge's Phenomenon")

analytically:  $X \Rightarrow X^T X$  also not invertible due to linear dependency of encoded categ. features perfect multicollinearity & dummy variable trap