

Solve non-linear problems by (approximate) reduction to linear ones

basic idea: if $Y = X \cdot \beta$ (linear model) is not good enough, add more features until it works, e.g. by the non-linear functions of X

e.g. $X = [X_1, X_2, X_3]$ (3 features originally)

$$\tilde{X} = [X_1, X_2, X_3, X_1^2, X_1 \cdot X_2, \log(X_3), \sqrt{X_2} \cdot \exp(X_3), \dots]$$

$\Rightarrow Y = \tilde{X} \cdot \tilde{\beta}$ depends non-linearly on X , despite being a linear fct.

improve by also transforming Y : $\tilde{Y} = \psi(Y, X)$ $\tilde{X} = \varphi(X)$

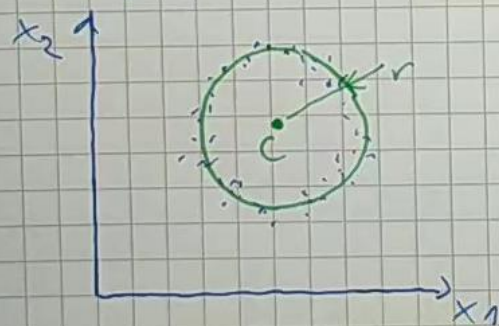
$$\text{solve } \tilde{Y} = \tilde{X} \cdot \tilde{\beta}$$

$$\text{predict: } Y = \psi^{-1}(\tilde{X} \cdot \tilde{\beta}, X)$$

ψ must be invertible in Y

Problem: coming up with good coordinate transforms $\tilde{Y} = \psi(Y, X)$, $\tilde{X} = \varphi(X)$ is difficult. Which functions to choose is a "hyperparameter"

Example where it works: circle fitting: 3 parameters: $\beta = [c_1, c_2, r]^T$



residual: distance of a point from the circle

$$\delta_i(\beta) = \sqrt{(x_{i1} - c_1)^2 + (x_{i2} - c_2)^2} - r = \epsilon \approx 0$$

$$\Rightarrow \boxed{\hat{\beta} = \arg \min_{\beta} \sum_{i=1}^N \delta_i(X_i)^2} \quad \text{no analytic solution}$$

\Rightarrow use Levenberg-Marquardt algorithm

alternative solution using coordinate transforms

Euclidean residual
↓

if point X_i is on the circle,
inside the circle
outside the circle

$$\sqrt{(x_{i1} - c_1)^2 + (x_{i2} - c_2)^2} - r = 0$$

$$< 0$$

$$> 0$$

likewise: X_i is on the circle
inside
outside

$$(x_{i1} - c_1)^2 + (x_{i2} - c_2)^2 - r^2 = 0$$

$$< 0$$

$$> 0$$

algebraic residual

⇒ find a coordinate transform that measure residuals with the second form (without square root)

$$\tilde{X}_i = [x_{i1}, x_{i2}, 1] = \varphi(x)$$

$$\tilde{Y}_i = x_{i1}^2 + x_{i2}^2 = \psi(x)$$

$$\tilde{\beta} = [2c_1, 2c_2, r^2 - c_1^2 - c_2^2]^T$$

$$= [2\beta_1, 2\beta_2, \beta_3^2 - \beta_1^2 - \beta_2^2]^T \quad (\beta = [c_1, c_2, r])$$

$$\Rightarrow \tilde{Y}_i - \tilde{X}_i \cdot \tilde{\beta} = (x_{i1} - c_1)^2 + (x_{i2} - c_2)^2 - r^2$$

$$\hat{\tilde{\beta}} = \underset{\tilde{\beta}}{\operatorname{argmin}} \sum_{i=1}^n (\tilde{Y}_i - \tilde{X}_i \cdot \tilde{\beta})^2$$

algebraic loss
solve as OLS

reconstructed original β :

$$\hat{c}_1 = \frac{\hat{\beta}_1}{2} \quad \hat{c}_2 = \frac{\hat{\beta}_2}{2} \quad r = \sqrt{\hat{\beta}_3^2 - \hat{c}_1^2 - \hat{c}_2^2}$$

when points are close to circle, both methods give similar results, algebraic fails for part of circles

Non-linear least squares in an augmented feature space

idea: $Y = f(X)$ is non-linear (hard to optimize), kind of coordinate transform $\tilde{X} = \phi(X)$ such that the relationship becomes (approximately) linear: $Y = \tilde{X} \cdot \beta \Rightarrow$ easy to optimize

[sometimes, it helps to transform Y as well: $\tilde{Y} = \tau(Y, X)$, can get desired response by inversion: $y = \tau^{-1}(\tilde{Y}, X)$, e.g. circle fit]

problem: finding a suitable transformation is hard

Q: Is there a systematic way of doing it? Universal recipe?

A: Yes. kernel ridge regression, kernel regression

intuition: if the optimization does not ~~not~~ access \tilde{X}_i directly, but only via scalar products $\tilde{X}_i \tilde{X}_j^T = G_{ij}$, then we can apply the kernel trick; compute G_{ij} via an analytic formula ("kernel function") without first computing \tilde{X}_i, \tilde{X}_j .

matrix of scalar products: Gram matrix $G = X X^T \in \mathbb{R}^{N \times N}$
Gramian $[G = \tilde{X} \tilde{X}^T]$

[note: don't confuse with scalar matrix $S = X^T X \in \mathbb{R}^{D \times D}$]

Normal ridge regression does not support the kernel trick:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} (Y - X\beta)^2 + \gamma \beta^T \beta \Rightarrow \hat{\beta} = \underbrace{(X^T X + \gamma I)^{-1} X^T Y}_{\text{regularized pseudo-inverse}}$$

because we need access to the feature vectors X_i

rewrite ~~big~~ ridge regression
dual solution vector $\hat{\alpha}$:

in terms of its dual optimization problem

$$\hat{\alpha} = \arg \max_{\alpha} - \alpha^T X X^T \alpha - \tau \mathbf{1}^T \alpha + 2 \alpha^T y$$

$$\hat{\alpha} = \underbrace{(X X^T + \tau \mathbb{I})}^{-1} y$$

$= G \Rightarrow$ access data only via Gram matrix
 $\hat{=}$ kernel can be applied

relation to the primal problem:

$$\hat{\beta} = X^T \hat{\alpha}, \quad \hat{y}_{\text{full}} = X_{\text{test}} \hat{\beta} = X_{\text{test}} X^T \hat{\alpha}$$

proof: \uparrow homework

g : scalar products
between test and
training features

$$g_i = X_{\text{test}} \cdot X_i^T$$

disadvantage of dual formulation: need to store training set also for
testing, so that $g = X_{\text{test}} X^T$ can be computed (not necessary
in primal ridge regression: once $\hat{\beta}$ is known, TS is no longer needed)

derivation of the dual problem

① introduce new variable v_i for the residuals: $v_i = y_i - X_i \beta$ for all i
 \Rightarrow primal problem $\beta = \arg \min_{\beta} r^T r + \tau \beta^T \beta$ s.t. $\forall_i: v_i = y_i - X_i \beta$

$$\Rightarrow \text{Lagrangian } (\beta, r, \gamma) = r^T r + \tau \beta^T \beta + \sum_{i=1}^n \gamma_i (y_i - X_i \beta - v_i)$$

now, it is mathematically more convenient, to factor out 2τ from the Lagrangian multipliers \Rightarrow introduce new multipliers α_i , s.t. $\beta_i = 2\tau \alpha_i$
 [always: $\beta_i \geq 0$, $\alpha_i \geq 0$] \Rightarrow modified Lagrangian

$$\text{Lagrangian}(\beta, r, \alpha) = r^T r + \tau \beta^T \beta + 2\tau \sum \alpha_i (y_i - x_i^T \beta - r_i)$$

② solve for β : $\frac{\partial \text{Lagrangian}}{\partial \beta} = 2\tau \beta - 2\tau \sum_{i=1}^n \alpha_i x_i^T \stackrel{!}{=} 0$

$$\hat{\beta} = \sum_i \alpha_i x_i^T = X X^T$$

\Rightarrow eliminate β from the Lagrangian

$$\begin{aligned} \text{Lagrangian}(r, \alpha) &= r^T r + \tau \alpha^T X X^T \alpha + 2\tau \alpha^T Y - 2\tau \alpha^T X^T X \alpha - 2\tau \alpha^T r \\ &= r^T r - \tau \alpha^T X X^T \alpha + 2\tau \alpha^T Y - 2\tau \alpha^T r \end{aligned}$$

③ solve for r : $\frac{\partial \text{Lagrangian}}{\partial r} = 2r^T - 2\tau \alpha^T \stackrel{!}{=} 0$ $r \neq 1/2 \alpha$ $r = \tau \alpha$

\Rightarrow eliminate r from Lagrangian:

$$\begin{aligned} \text{Lagrangian}(\alpha) &= \tau^2 \alpha^T \alpha - \tau \alpha^T X X^T \alpha + 2\tau \alpha^T Y - 2\tau^2 \alpha^T \alpha \\ &= -\tau \alpha^T X X^T \alpha + 2\tau \alpha^T Y - \tau^2 \alpha^T \alpha \end{aligned}$$

\Rightarrow dual optimization problem $\hat{\alpha} = \arg \max_{\alpha} -\alpha^T X X^T \alpha + 2\alpha^T Y - \tau \alpha^T \alpha$

④ Solve dual for α : \downarrow Lagrang. $= -2 X X^T \alpha - 2 \epsilon \alpha + 2 Y \stackrel{!}{=} 0$

$$\boxed{2 = (X X^T + \epsilon I)^{-1} Y}$$

$$= (G + \epsilon I)^{-1} Y \Rightarrow \text{can apply kernel trick}$$

Example for a coordinate transform where the kernel trick works

given: 2-D feature space $X_i = [X_{i1}, X_{i2}]$

augmented feature space: $\tilde{X}_i = \varphi(X_i) \in \mathbb{R}^6$

$$= [\sqrt{2} X_{i1}, \sqrt{2} X_{i2}, X_{i1}^2, X_{i2}^2, \sqrt{2} X_{i1} \cdot X_{i2}, 1]$$

dot product: $\tilde{X}_i \cdot \tilde{X}_{i'} = 2 X_{i1} \cdot X_{i'1} + 2 X_{i2} \cdot X_{i'2} + X_{i1}^2 \cdot X_{i'1}^2 + X_{i2}^2 \cdot X_{i'2}^2 + 2 X_{i1} \cdot X_{i2} \cdot X_{i'1} \cdot X_{i'2} + 1$

$$= ((X_i \cdot X_{i'}^T) + 1)^2$$

Gram matrix in augmented space: $G = \tilde{X} \cdot \tilde{X}^T$

with $G_{i,i'} = (X_i \cdot X_{i'}^T + 1)^2 = k(X_i, X_{i'})$

\Rightarrow can calculate $G_{i,i'} = \tilde{X}_i \cdot \tilde{X}_{i'}^T$ without computing \tilde{X}_i and $\tilde{X}_{i'}$

a formula with this property is called kernel function: $G_{i,i'} = k(X_i, X_{i'})$

Mercer condition: a formula $k(X_i, X_{i'})$ is a valid kernel function (i.e., is identical

to a scalar product in some augmented feature space), when for all square-integrable

function $f(x)$ it holds: $\int \int k(X_i, X_{i'}) f(X_i) \cdot f(X_{i'}) dX_i dX_{i'} \geq 0$

Intuition: for any FS $\{X_i, X_{i'}\}$, the Gram matrix $G_{i,i'} = k(X_i, X_{i'})$ is positive semi-definite

⇒ when Mercer's condition is fulfilled, the inverse $(K + \tau I)^{-1}$ exists
for every training set ($\hat{z} = (K + \tau I)^{-1} y$ exists)

• most popular kernel functions

- polynomial kernels: $K(x_i, x_{i'}) = (x_i^T x_{i'} + 1)^p$ $p \geq 1$
(the example from previous page used $p=2$)

- squared exponential kernel $K(x_i, x_{i'}) = \exp\left(-\frac{\|x_i - x_{i'}\|^2}{2h^2}\right)$

h : band width of kernel \Rightarrow how the distance $\|x_i - x_{i'}\|$ is scaled

- p or h (and τ from ridge regression) are hyper-parameters that the user must adjust (e.g. cross-validation)

- long list of permissible kernel function by asking Google

- more general treatment: Gaussian processes in Advanced Machine Learning

• heuristic approximation: kernel regression (Nahavaya-Watson regression)

$$\hat{y}_{new} = \hat{f}(x_{new}) = \frac{\sum_{i=1}^N y_i \cdot g_i}{\sum_{i=1}^N g_i} \quad \begin{matrix} g_i = K(x_{new}, x_i) \\ \text{weighted sum} \end{matrix}$$

complexity: $O(N \cdot D)$

$g_i' = \frac{g_i}{\sum_{i=1}^N g_i}$
weight normalization

• in practice (\hat{f} squared-exponential kernel), $g_i \approx 0$ when $\|x_i - x_{new}\| > 3h$

\Rightarrow only the near-neighbors of x_{new} contribute to the sum

\Rightarrow efficient algorithms to sum only over near neighbors (instead of all N training points)
e.g. "fast Gaussian transform"

- prediction of new distance with dual ridge regression.

- linear case: $\hat{y} = X_{\text{new}} X^T \hat{\alpha} \Rightarrow$ pre-compute $\hat{\beta}^T = X^T \alpha$
 $= X_{\text{new}} \hat{\beta} \Rightarrow$ do not need to store X (TS)

- non-linear case (with kernel trick): $\hat{y} = \underbrace{K(X_{\text{new}}, X)}_{= g} \hat{\alpha}$
 with $g_i = K(X_{\text{new}}, x_i)$

$\Rightarrow g$ must be computed again for every x_{new} , pre-computing $X^T \hat{\alpha}$ does not help, because $K(\cdot, \cdot)$ is non-linear, and we cannot factor out $X^T \hat{\alpha}$ as in the linear case \Rightarrow need to store training set

- when kernel trick was invented in the 1990s, explicitly computing high-dimensional \tilde{X} was impractical (too expensive in time and memory)

\Rightarrow kernel avoids this

today, we have GPUs and lots of RAM \Rightarrow we can explicitly compute

complex $\tilde{X} = \phi(X)$, for example many neural nets do this:

network with L layers: - first $L-1$ layers compute $\tilde{X} = \phi(X)$ \leftarrow
 - last layer computes $\hat{y} = \tilde{X} \cdot \beta$

per ultimate neuron activations form an augmented feature space: "latent representation"

- kernel trick can also be used to make linear classifiers (LDA, LR) non-linear \Rightarrow support vector machine (SVM) \Rightarrow later