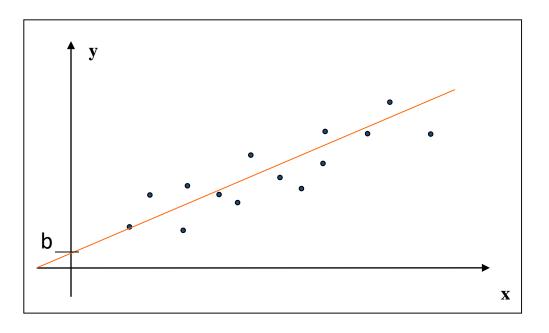


From Linear Regression to Weighted Nonlinear Regression



Linear regression searches a linear mapping between input x and output y, parametrized by the slope vector w and intercept b.

$$y = f(x; w, b) = w^T x + b$$





Linear regression searches a linear mapping between input x and output y, parametrized by the slope vector w and intercept b.

$$y = f(x; w, b) = w^T x + b$$

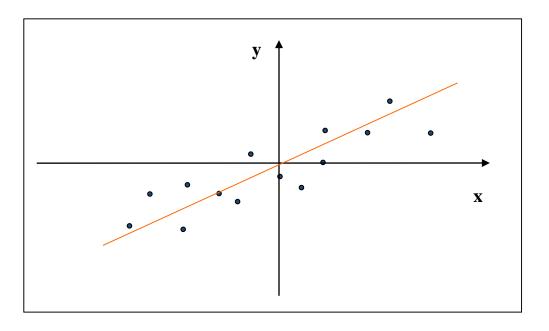
One can omit the intercept by centering the data:

$$y' = y - \overline{y}$$
 and $x' = x - \overline{x}$, \overline{x} , \overline{y} : mean on x and y
 $y' = w^T x' + b'$
with $b' = b + w^T \overline{x} - \overline{y}$
Least-square estimate of $(b')^* = \overline{y}' - w^T \overline{x}' = 0$
 $\Rightarrow y' = w^T x'$.



Linear regression searches a linear mapping between input *x* and output *y*, parametrized by the slope vector *w*.

$$y = f(x; w) = w^T x$$





Pair of M training points $X = [x^1 \ x^2 \ ... \ x^M]$ and $y = [y^1 \ y^2 \ ... \ y^M]^T$ $x^i \in \mathbb{R}^N, y^i \in \mathbb{R}.$

Find the optimal parameter w through least-square regression:

$$w^* = \min_{w} \left(\sum_{i=1}^{M} \frac{1}{2} (w^T x^i - y^i)^2 \right)$$

Find the analytical solution through partial differentiation:

$$w^* = \left(XX^T\right)^{-1}Xy$$

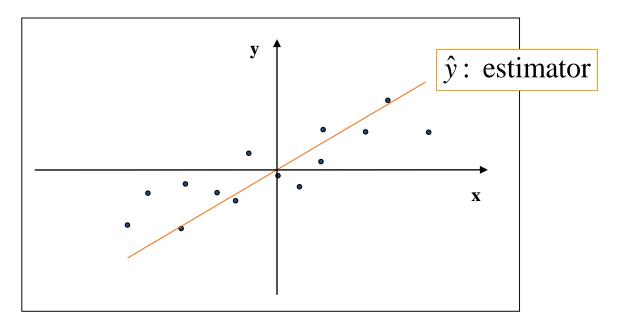


Weighted Linear Regression

Regression through weighted Least Square

$$w^* = \min_{w} \left(\sum_{i=1}^{M} \frac{1}{2} \beta_i \left(w^T x^i - y^i \right)^2 \right), \quad \beta_i \in \mathbb{R} \quad \beta_1 = \beta_2 ... = \beta_M$$

$$\Rightarrow \text{Standard linear regression}$$



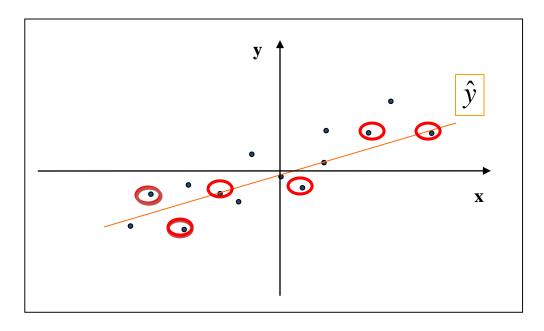
All points have equal weight.



Weighted Linear Regression

Regression through weighted Least Square

$$w^* = \min_{w} \left(\sum_{i=1}^{M} \frac{1}{2} \beta_i \left(w^T x^i - y^i \right)^2 \right), \quad \beta_i \in \mathbb{R}$$



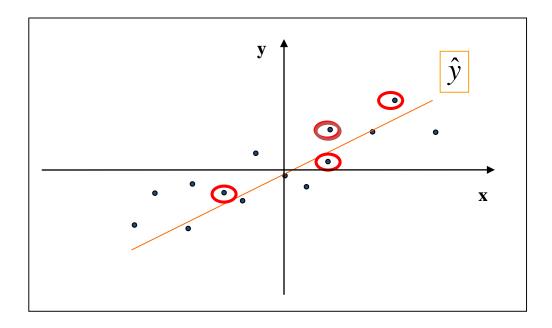
Points in red have large weights.



Weighted Linear Regression

Regression through weighted Least Square

$$w^* = \min_{w} \left(\sum_{i=1}^{M} \frac{1}{2} \beta_i \left(w^T x^i - y^i \right)^2 \right), \quad \beta_i \in \mathbb{R}$$



Points in red have large weights.



Computing the optimal regressor

We start from the loss function: $\sum_{i=1}^{M} \frac{1}{2} \beta_i \left(w^T x^i - y^i \right)^2$

We set B a diagonal matrix with entries
$$\beta_i$$
, $B = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \\ \beta_M \end{bmatrix}$

Change of variable: $Z = XB^{1/2}$ and $v = B^{1/2}y \implies \text{Loss: } w^TZ - v$

Minimizing for loss, one gets a closed-form best estimator for w:

$$\hat{y} = w^T x = \left(Z Z^T \right)^{-1} Z v x$$

$$w^* = \left(ZZ^T\right)^{-1}Zv$$

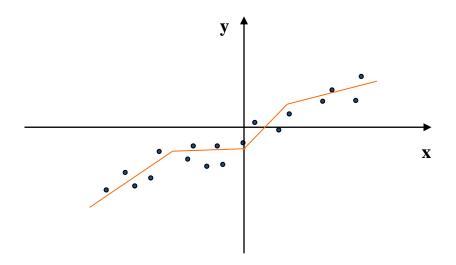


Limitations of Linear Regression

$$w^* = \min_{w} \left(\sum_{i=1}^{M} \frac{1}{2} \beta_i \left(w^T x^i - y^i \right)^2 \right), \quad \beta_i \in \mathbb{R} \text{ constant}$$

Assumes that a single linear dependency applies everywhere.

Not true for data sets with local dependencies.



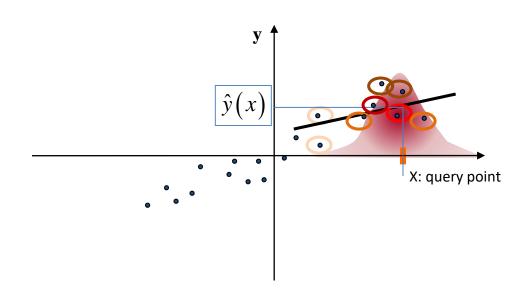
Need a regression method to estimate local linear dependencies



Estimate is determined through local influence of each group of datapoints

$$\widehat{y}(x) = \sum_{i=1}^{M} \beta_i(x) y^i / \sum_{j=1}^{M} \beta_j(x) \qquad \beta_i(x) \in \mathbb{R} : \text{ weights function of } x$$

$$\beta_i(x) = K(d(x^i, x)), \text{ with } K(d(x^i, x)) = e^{-d(x^i, x)}, d(x^i, x): \text{norm-2}$$

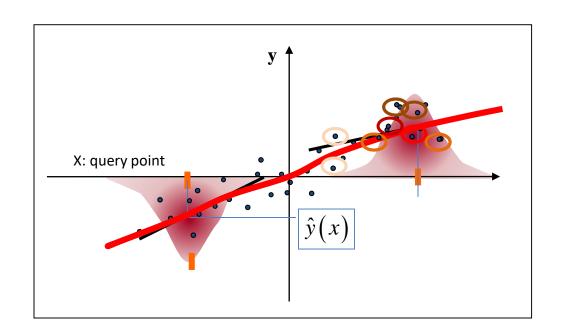




Estimate is determined through local influence of each group of datapoints

$$\widehat{y}(x) = \sum_{i=1}^{M} \beta_i(x) y^i / \sum_{j=1}^{M} \beta_j(x) \qquad \beta_i(x) \in \mathbb{R}: \text{ weights function of } x$$

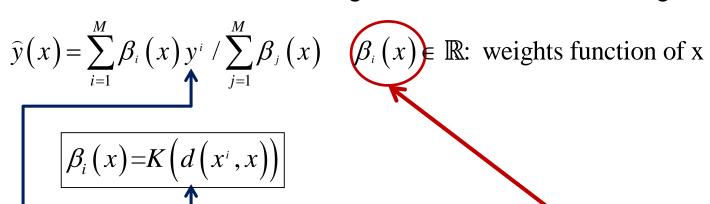
$$\beta_i(x) = K(d(x^i, x))$$
, with $K(d(x^i, x)) = e^{-d(x^i, x)}$, $d(x^i, x)$:norm-2



Generates a smooth function y(x)



Estimate is determined through local influence of each group of datapoints



Model-free regression! No longer explicit model of the form $y = w^T x$

Regression computed at each query point. Depends on training points.



Estimate is determined through local influence of each group of datapoints

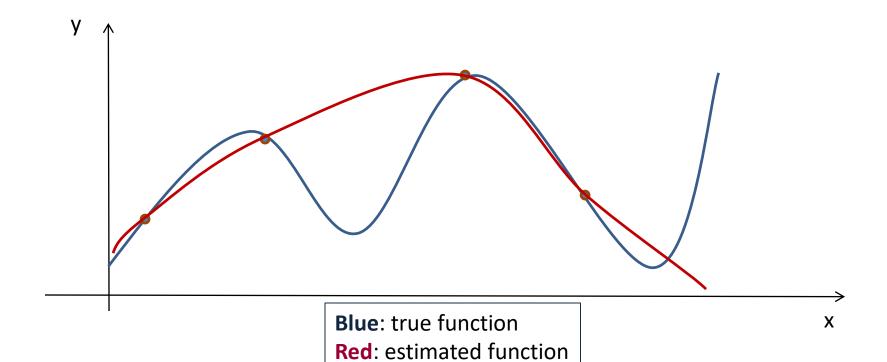
$$\widehat{y}(x) = \sum_{i=1}^{M} \beta_{i}(x) \quad y^{i} \sum_{j=1}^{M} \beta_{j}(x) \quad \beta_{i}(x) \in \mathbb{R}: \text{ weights function of } x$$

$$Which training points?$$

$$Which kernel?$$



Good prediction depends on the choice of datapoints.

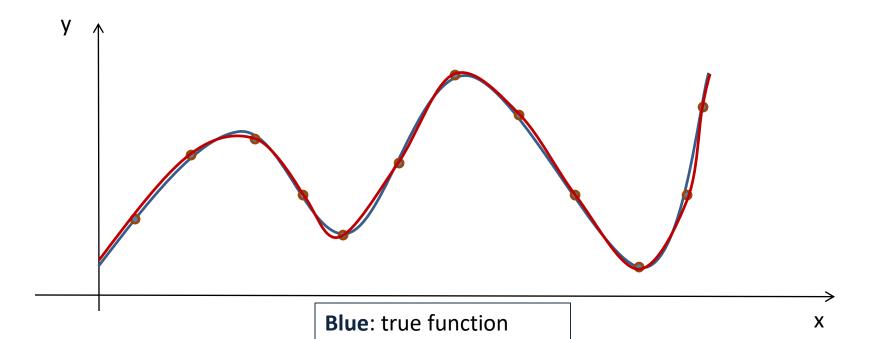




Good prediction depends on the choice of datapoints.

The more datapoints, the better the fit.

Computational costs increase dramatically with number of datapoints



Red: estimated function

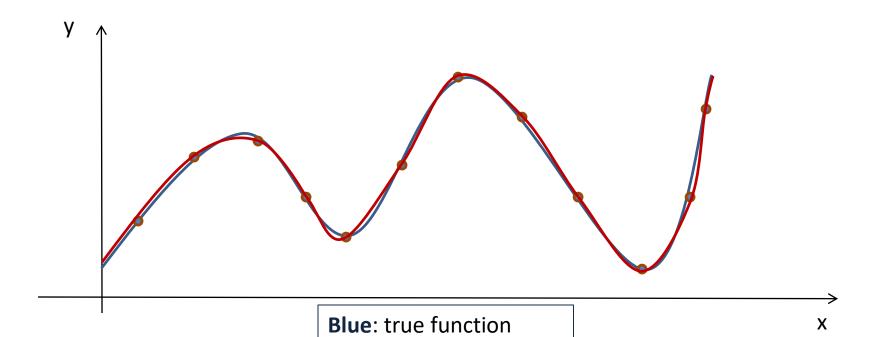
16



Several methods in ML for performing non-linear regression.

Differ in the objective function, in the amount of parameters.

K-nearest neighbors (KNN) uses all datapoints (model-free)



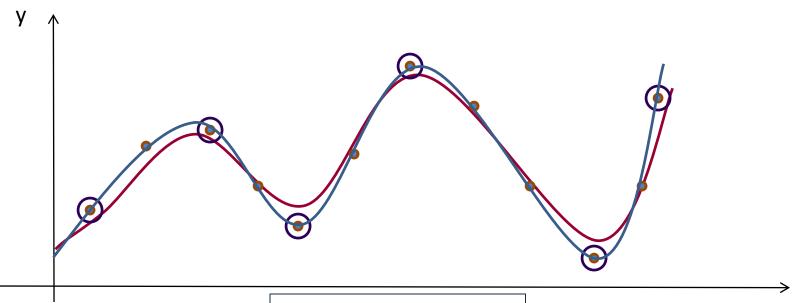
Red: estimated function



Several methods in ML for performing non-linear regression.

Differ in the objective function, in the amount of parameters.

K-nearest neighbors (KNN) uses all datapoints (model-free)
Support Vector Regression (SVR) picks a subset of datapoints (support vectors)



Blue: true function

Red: estimated function

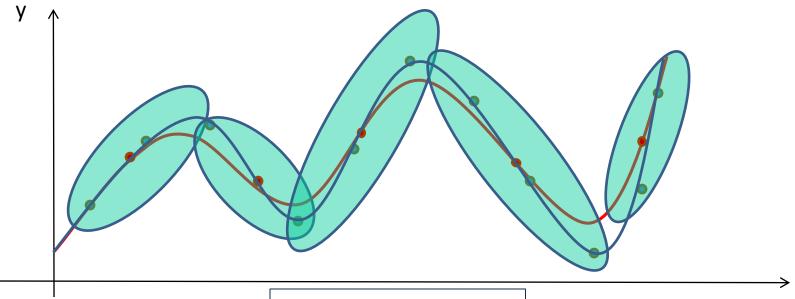
Χ



Several methods in ML for performing non-linear regression.

Differ in the objective function, in the amount of parameters.

K-nearest neighbors (KNN) uses all datapoints (model-free)
Support Vector Regression (SVR) picks a subset of datapoints (support vectors)
Gaussian Mixture Regression (GMR) generates a new set of datapoints (centers of Gaussian functions)

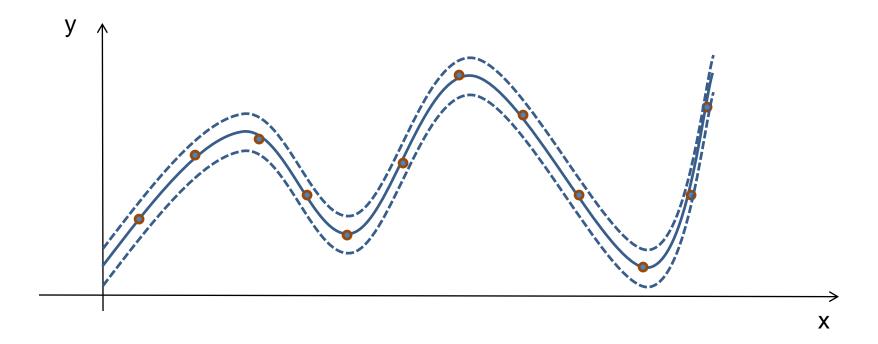


Blue: true function

Red: estimated function



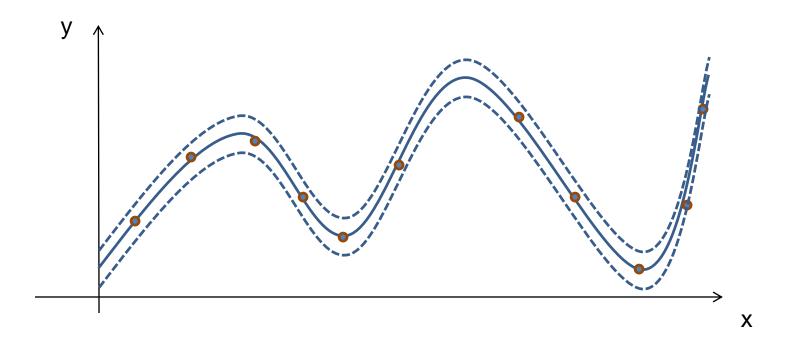
Estimate of the noise is important to measure goodness of fit.





Estimate of the noise is important to measure goodness of fit.

Support Vector Regression (SVR) assumes an estimate of the noise model (ε -tube) and then compute f directly within a noise-tolerance.





Estimate of the noise is important to measure goodness of fit.

Gaussian Mixture Regression (GMR) builds a local estimate of the noise model through the variance of the system.

