

Recall, we assume for regression

$$y = f(x) + \varepsilon$$

want to determine \hat{f} so that $\hat{y} = \hat{f}(x) \approx y$.

Solu! minimize expected loss (Risk)

$$\hat{f} = \arg \min_f E[(y - f(x))^2]$$

claim: $\hat{f}(x) = E[y|x]$

So we model $E[y|x]$.

Ways to do this:

① Linear Regression

$$\hat{f}(x) = x^T \hat{\beta}$$

② Polynomial Regression $x \in \mathbb{R}$ ($p=1$)

$$\hat{f}(x) = \hat{\beta}_0 + \hat{\beta}_1 x + \hat{\beta}_2 x^2 + \hat{\beta}_3 x^3 + \dots$$

Global
assumption
about
form
of
 $E[y|x]$

Local
assumption

③ relax this global (linear) assumption

local
assumption

(o) relax this global (linear) assumption
e.g. KNN regression.

→ non-parametric regression

model:

$$E(y|x) = m(x)$$

↑ some function w/ no
particular predefined
form.

KNN Regression

$$\hat{f}(x) = \frac{1}{K} \sum_{n: \underline{x}_n \in N_K(x)} y_n$$

if we let

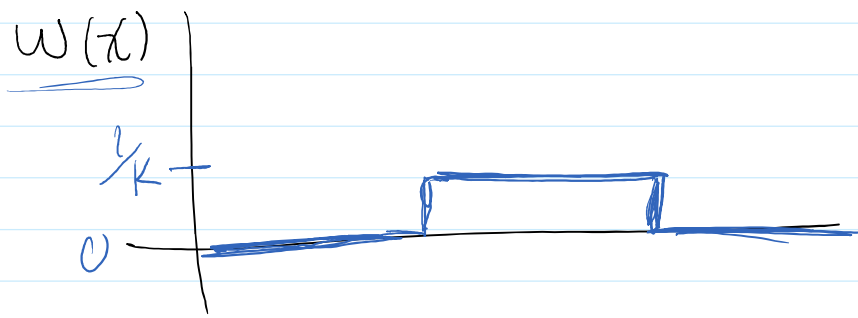
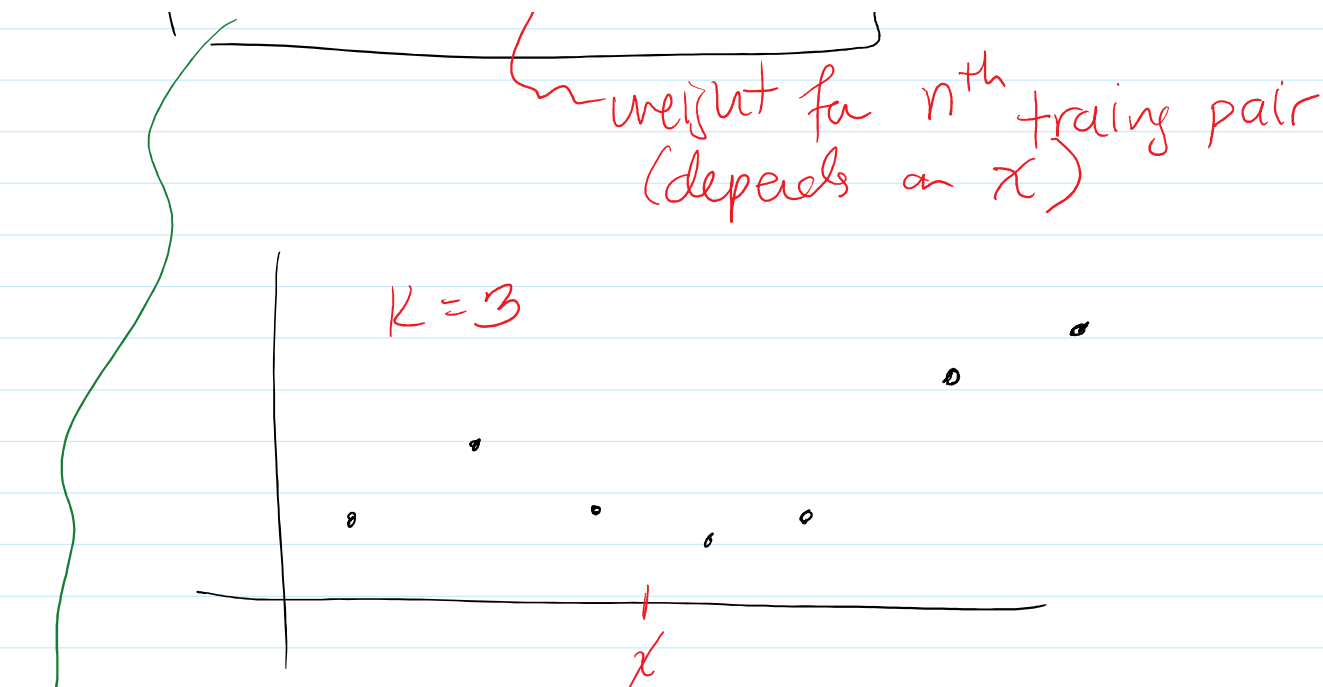
$$\begin{aligned} \omega_n(x) &= \frac{1}{K} \mathbb{1}(x_n \in N_K(x)) \\ &= \begin{cases} 1/K, & x_n \in N_K(x) \\ 0, & \text{else} \end{cases} \end{aligned}$$

then for KNN

$$\hat{f}(x) = \sum_{n=1}^N \omega_n(x) y_n$$

weighted average of
training y 's

weight of n^{th} ... pair

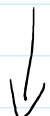


KNN may be called a linear smoother

Notice! w_n s are discontinuous.

Generalization is to use "nicer" weighting function

w_n



continuous
differentiable

leads to Kernel Smoothing.

Defn: a Kernel is a function K

① $K: \mathbb{R}^P \rightarrow \mathbb{R}$

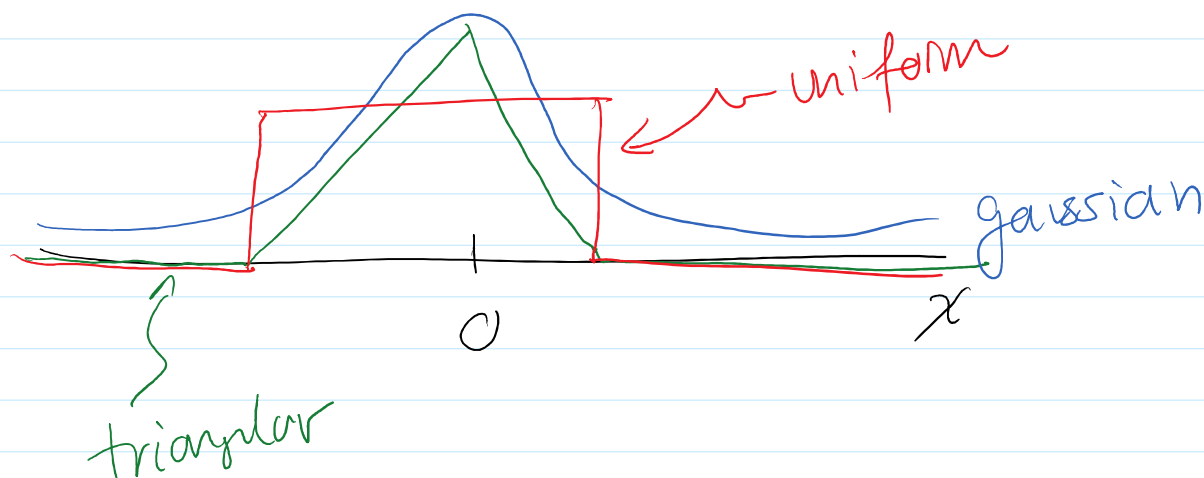
② $K(x) \geq 0$

typically also assume

③ $\int_{\mathbb{R}^P} K(x) dx = 1$ (PDF of a P -variate RV)

④ $K(x) = K(-x)$ (symmetric)

Ex.



We can use kernel, to define a measure of "closeness":

$$K(x-u) \approx \text{closeness between } x \text{ and } u$$

$K(x-y) \approx$ closeness between x and y

$p=1$ and $K(\cdot) = \frac{1}{|\cdot|}$ then $K(x-y) = \frac{1}{|x-y|}$.

Ex. $K(\cdot) = \exp(-\cdot^2)$

$$K(x-y) = \exp(-(x-y)^2)$$

linear smoother:

Kernel Smoothing

$$\hat{f}(x) = \sum_{n=1}^N w_n(x) y_n$$

where $w_n(x) \propto \underbrace{K(x-x_n)}_{\text{closeness of } x_n \text{ to } x}$

Typically req that

$$\sum_{n=1}^N w_n(x) = 1$$

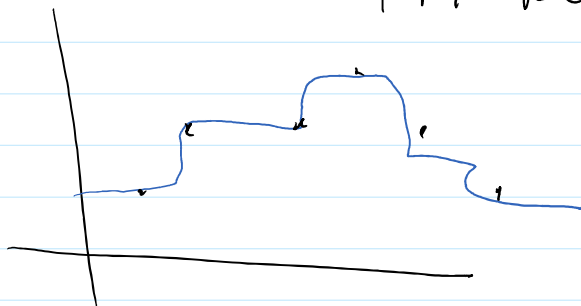
i.e.

$$w_n(x) = \frac{K(x-x_n)}{\sum_{i=1}^N K(x-x_i)}$$

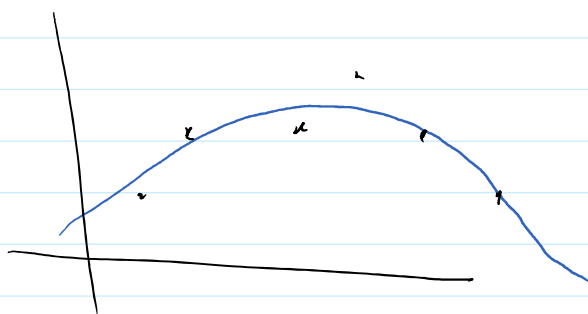
→ Nadaraya-Watson (NW)

Kernel Regression estimator.

KNN Reg



Kernel Reg



Generalize a Bit more

For weights w_n then the weighted avg of y_n s

$$\rightarrow \bar{y}_w = \frac{1}{\sum_n w_n} \sum_{n=1}^N w_n y_n$$

if $w_n = 1/N \forall n$ we get typical avg.

Model : $Y \approx \alpha$

$$\hat{\alpha} = \underset{\alpha}{\operatorname{argmin}} \sum_n (y_n - \alpha)^2 = \bar{y}$$

$$\rightarrow \hat{\alpha} = \underset{\alpha}{\operatorname{argmin}} \sum_n \underline{w_n} (y_n - \underline{\alpha})^2 = \bar{y}_w$$

Model: $Y \approx \alpha + \beta x$

$$\hat{\alpha}, \hat{\beta} = \underset{\alpha, \beta}{\operatorname{argmin}} \sum_n (y_n - \alpha - \beta x_n)^2 = \text{Least squares reg. ests}$$

$$\hat{\alpha}, \hat{\beta} = \underset{\alpha, \beta}{\operatorname{argmin}} \sum_n w_n (y_n - \alpha - \beta x_n)^2 = \text{weighted LS reg. ests.}$$

↑ can depend on x_n s

Model: $Y = X\beta$

$$\hat{\beta}^{(ols)} = \underset{\beta}{\operatorname{argmin}} \|Y - X\beta\|^2 = (X^T X)^{-1} X^T Y$$

$$\hat{\beta}^{(wls)} = \underset{\beta}{\operatorname{argmin}} \| \overset{\text{diag}(w_n)}{\leftarrow} W(Y - X\beta) \|^2 = (X^T W X)^{-1} X^T W Y$$

Punchline: $x \in \mathbb{R}$ ($p=1$)

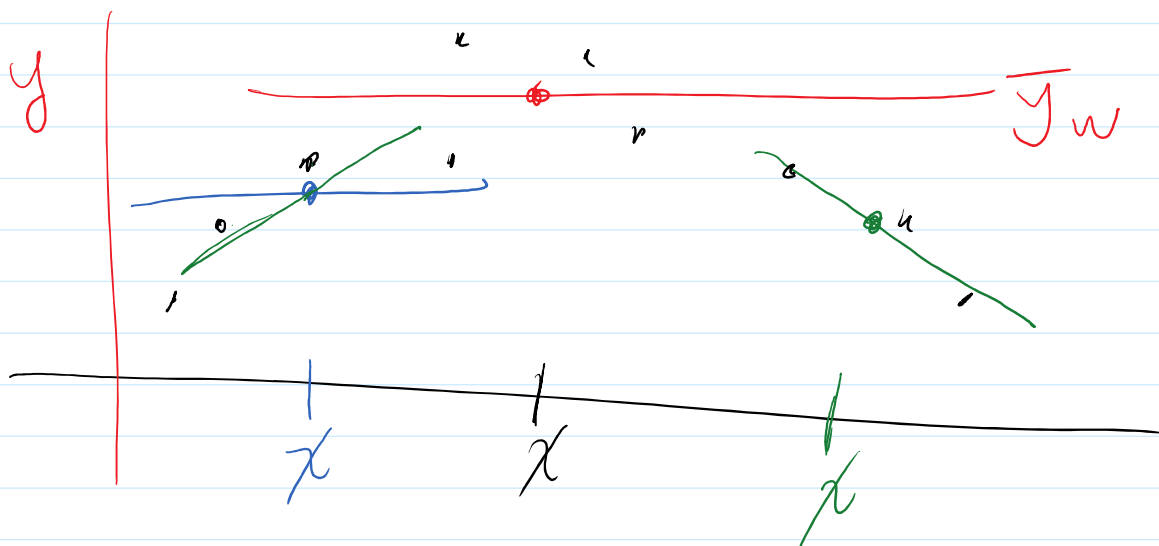
$$\hat{f}(x) = \underset{\theta}{\operatorname{argmin}} \sum_{n=1}^N \underbrace{K(x - x_n)}_{w_n} (y_n - \theta)^2$$

then

$$\hat{f}(x) = \bar{y}_w \quad w/ \quad w_n \propto K(x - x_n)$$

= NW Kernel Regression est.
w/ $w_n(x) \propto K(x - x_n)$.

NW K-Reg. est. solves a simple weighted regression model where weights depend on x so that $w_n(x) \propto K(x - x_n)$



Local Polynomial Regression

Simple example

$$\hat{f}(x) = \hat{\alpha}(x) + \hat{\beta}(x)x$$

$$\hat{\alpha}(x), \hat{\beta}(x) = \underset{\alpha, \beta}{\operatorname{argmin}} \sum_{n=1}^N \overbrace{K(x-x_n)}^{w_n(x)} (y_n - \alpha - \beta x_n)^2$$

weighted regression ests
w/ $w_n(x) \propto K(x-x_n)$

Fully $P=1$ Local Poly

$$\{\hat{\beta}_i(x)\}_i = \underset{\beta_1, \dots, \beta_d}{\operatorname{argmin}} \sum_{n=1}^N K(x-x_n) \left(y_n - \sum_{j=1}^d \beta_j x_n^j \right)^2$$

$$\sum_j \{p_j^2 S_j\} \quad n=1$$

weighted poly regression

ultimately: at some x

$$W = \text{diag}(W_n(x)) = \text{diag}(K(x - x_n))$$

$$X = \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 & \dots \\ 1 & x_2 & x_2^2 & x_2^3 & \dots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

ad $\hat{\beta} = (X^T W X)^{-1} X^T W y$

ad so

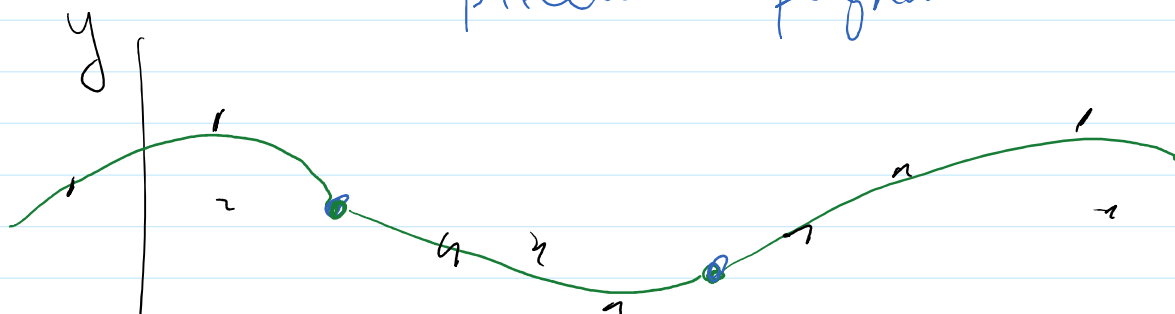
$$\hat{y} = \hat{f}(x) = x^T \hat{\beta}$$

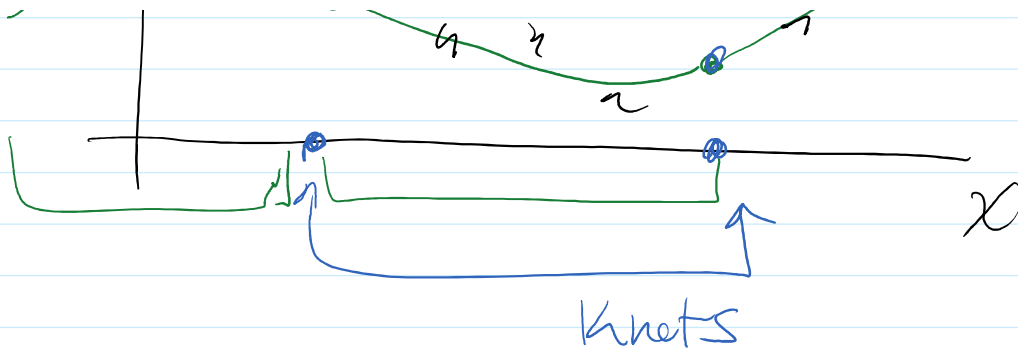
linear smoother.

Splines

Kernel Reg/Local Poly: local methods

Splines: Fit a better global model.
piecewise-polynomials





Spline: k -degree polynomial w/ some
smoothness condition at break-points
