3) Mean Shift

a) $x_j^{t+1} = x_j^t + \alpha_j^t \frac{2}{n} \sum\limits_{i: \|x_i - x_j^t\| < 1} (x_i - x_j^t)$

Goal: update to be equivalent to moving to the local mean of
points within the kernel support

Idea: set $\alpha_j^t$ so that update vector moves $x_j^t$ directly to
the mean of nearby points.

$\nu_j = \sum\limits_{i: \|x_i - x_j^t\| < 1} 1$ ; local mean $= \frac{1}{\nu_j} \sum\limits_{i: \|x_i - x_j^t\| < 1} x_i$

$\longmapsto$ local mean $- x_t^j = \frac{1}{\nu_j} \sum\limits_{i: \|x_i - x_j^t\| < 1} (x_i - x_{t_j}^t)$

$x_{t+1}^j = x_t^j + \alpha_j^t \frac{2}{n} \nu_j \left( \frac{1}{\nu_j} \sum\limits_{i: \|x_i - x_j^t\| < 1} (x_i - x_j^t) \right)$

$\phantom{x_{t+1}^j} = x_j^t + \alpha_j^t \cdot \frac{2\nu_j}{n} (\text{local mean} - x_t^j)$.

$\Rightarrow \alpha_j^t = \frac{n}{2\nu_j}$

sensible because:

• learning rate adapts to the local density of points
   many nearby points → small stepsize and vice versa

• ensures that $x_j^t$ is moved directly to the mean of nearby
   points ⇒ convergence to local modes.

# 5) Linear Regression : Heteroscedastic Noise

$$y_n = \beta^T x_n + \varepsilon_n \qquad E[\varepsilon_n] = 0 \qquad Var[\varepsilon_n] = \sigma_n^2$$

minimize the weighted sum of squares of the residuals

$$J(\beta) = \sum_{n=1}^{N} (y_n - \beta^T x_n)^2 \longrightarrow \sum_{n=1}^{N} (y_n - \beta^T x_n)^2 / \sigma_n^2$$

observations with lower variance contribute more to the determination of the coefficients $\beta$ that the ones with higher variance. By minimizing $J(\beta)$ we find the best estimate.

$$\hat{\beta} = (X^T W X)^{-1} X^T W y \qquad \text{with } W \text{ diagonal } W_{nn} = \frac{1}{\sigma_n^2}$$

$y$ observed values $x$ matrix of predictors

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

$$= (X^T W X)^{-1} X^T W \underbrace{E[y]}_{= (X^T W X)^{-1} X^T W (X\beta)}$$

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

$$Cov(Ay) = A \, Cov(y) \, A^T \quad ; \quad Cov(y) = Cov(x\beta + \varepsilon) = Cov(\varepsilon) = diag(\sigma_1^2, \sigma_2^2, \ldots \sigma_N^2)$$

$$Cov(y) = \sigma^2 I$$

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

$$= (X^T W X)^{-1} X^T W W X (X^T W X)^{-1} \sigma^2$$

$$= (X^T W X)^{-1} X^T W X (X^T W X)^{-1} \sigma^2 = \sigma^2 (X^T W X)^{-1}$$

$\sigma^2$ is a scaling factor, this accounts for the different variances.