# Prediction

Setup: Model relating $y$ to $(1, x_1, \ldots, x_p)^T = \underline{x}$:

$$Y = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p + \varepsilon = \beta^T \underline{x} + \varepsilon$$

and OLS estimates $\hat{\beta}$ (plus something for $\hat{\sigma}^2$),

there are two new quantities we want to predict at a new set of features $\underline{x}_* = (1, x_{1*}, \ldots, x_{*p})^T$:

- Average response / mean value $\qquad \beta^T \underline{x}_* \qquad \color{red}{[fit]}$

- New obs $\qquad \beta^T \underline{x}_* + \varepsilon_* \quad \color{red}{[prediction]}$

Both cases use same point predictor:

$$\hat{\beta}^T \underline{x}_* = \underline{x}_*^T \hat{\beta} = \underline{x}_* (X^T X)^{-1} X^T y$$

but, uncertainty depends on the case!

Recall standard error is an estimate of standard deviation

$$SE(\cdot) = \sqrt{\widehat{Var}(\cdot)}$$

and

$$Var(\underline{x}^T \hat{\beta}) = Cov(\underline{x}^T \hat{\beta}, \underline{x}^T \hat{\beta}) = \underline{x}^T (Var \hat{\beta}) \underline{x}$$

So,

$$\widehat{Var}(\underline{x}_*^T \hat{\beta}) = \hat{\sigma}^2 \underline{x}_*^T (X^T X)^{-1} \underline{x}_* \qquad \text{\textcolor{red}{[fit]}}$$

$$\widehat{Var}(\underline{x}_*^T \hat{\beta} + \varepsilon_*) = \hat{\sigma}^2 \underline{x}_*^T (X^T X)^{-1} \underline{x}_* + \hat{\sigma}^2 \qquad \text{\textcolor{red}{[prediction]}}$$
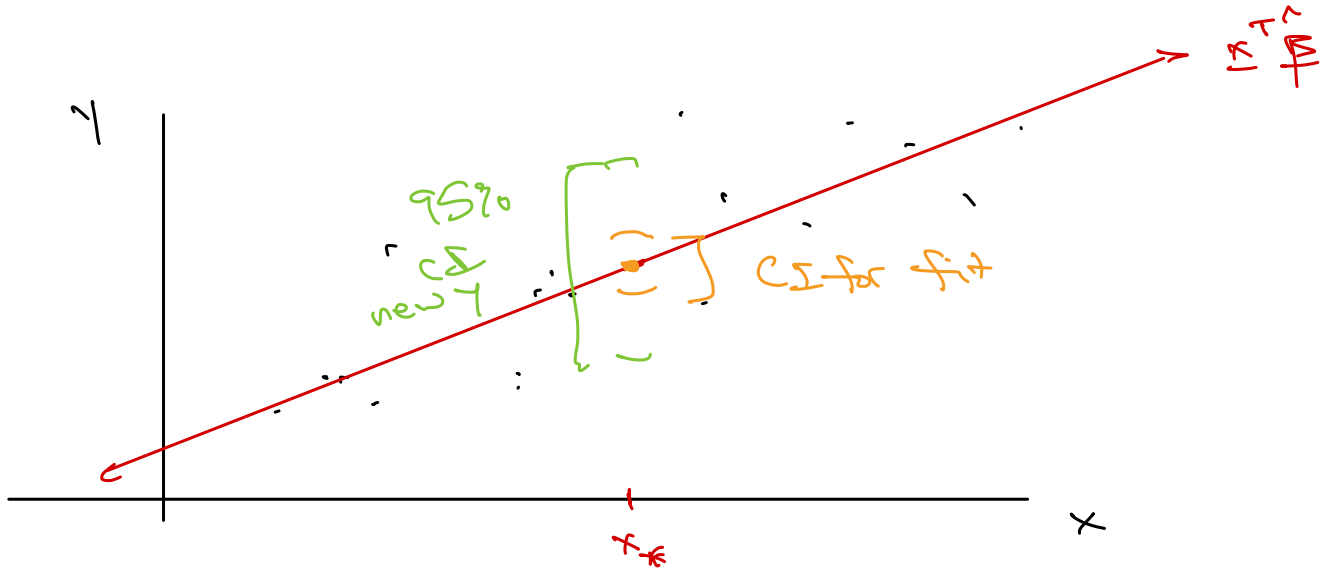
Thus 95% approx conf. intervals are:

fit:
$$x_*^T \hat{\beta} \pm 1.96 \, SE\left(x_*^T \hat{\beta}\right)$$

prediction:
$$x_*^T \hat{\beta} \pm 1.96 \, SE\left(x_*^T \hat{\beta} + \Sigma_*\right)$$

## 4.5 Diagnostics

Have model

$$\underline{Y} = X \underline{\beta} + \underline{\varepsilon}$$

+ estimators $\hat{\underline{\beta}}, \hat{\sigma}^2$. Is the model any good, and are $\{\varepsilon_i\}$ approximately normal?
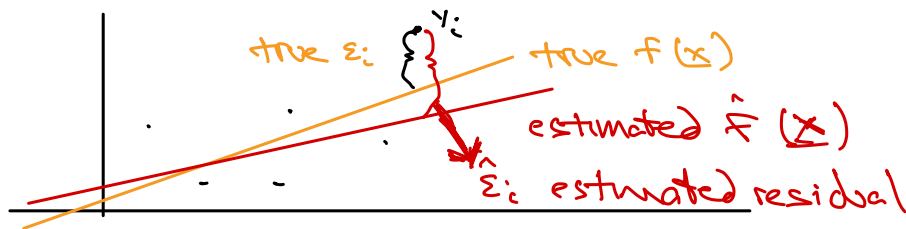
## DEF

Fitted values $\qquad \hat{Y}_i = x_i^T \hat{\underline{\beta}}$, $\qquad \hat{\underline{Y}} = X \hat{\underline{\beta}}$ $\qquad [\neq \underline{Y} = X\underline{\beta} + \underline{\varepsilon}]$
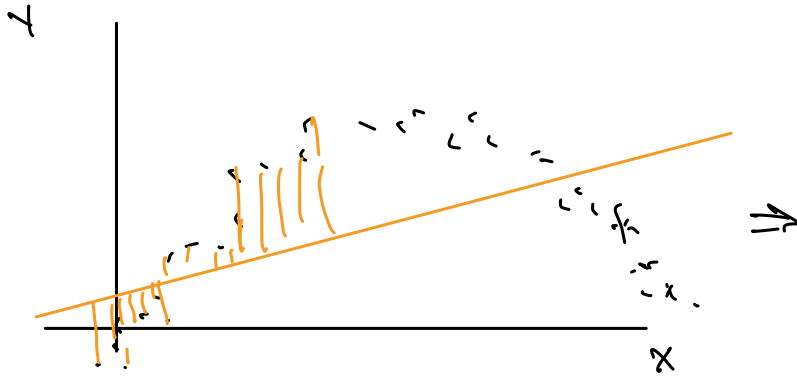
Estimated residuals $\qquad \hat{\underline{\varepsilon}} = \underline{Y} - X \hat{\underline{\beta}}$ $\qquad [\neq \underline{\varepsilon} = \underline{Y} - X\underline{\beta}]$



true $\varepsilon_i$ $\qquad$ true $f(x)$

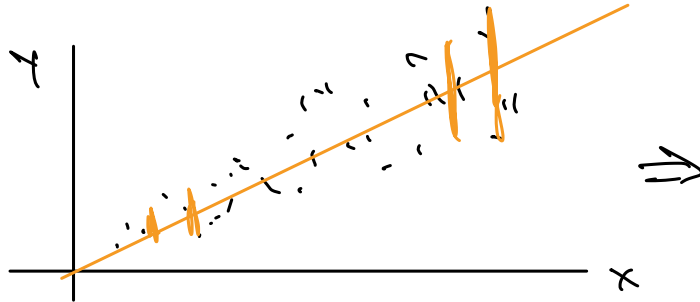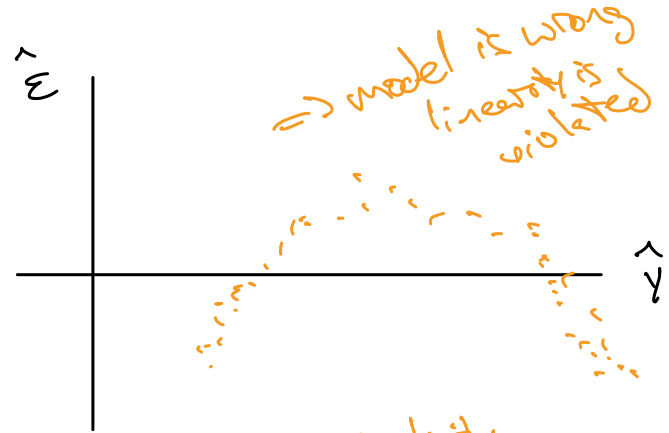estimated $\hat{f}(x)$

$\hat{\varepsilon}_i$ estimated residual

If assumptions are correct, $\hat{\varepsilon}$ should have no structure.
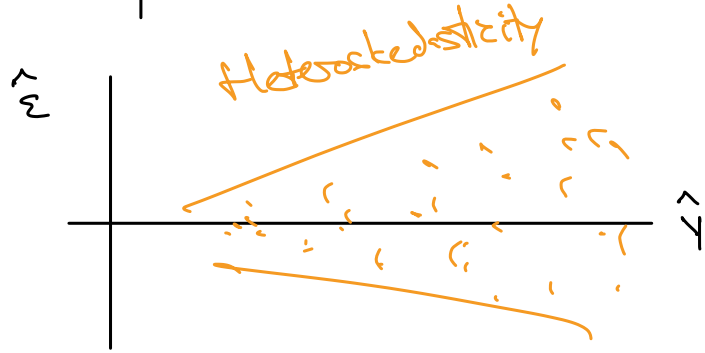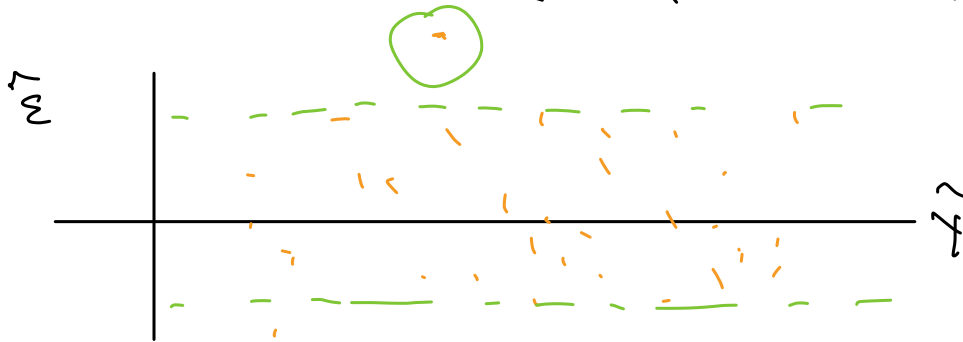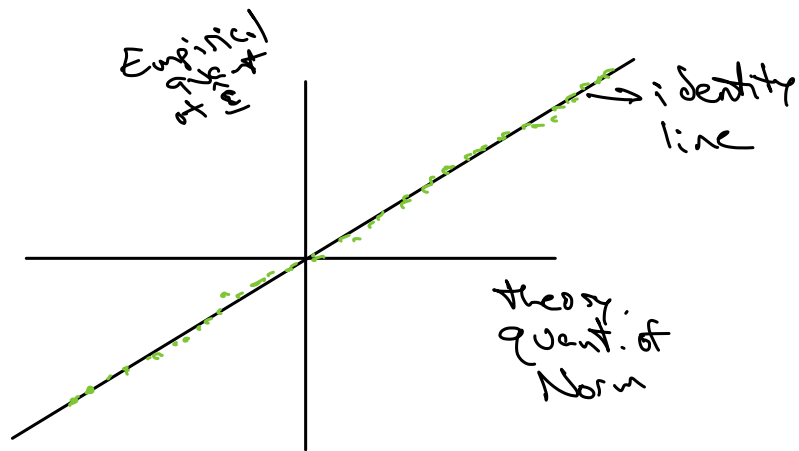
① Fitted values vs. estimated! residuals plot

② Unusual values of $|\hat{\varepsilon}|$ may be evidence of <u>outliers</u>



Outliers do not necessarily affect $\hat{\beta}$, but do inflate $\hat{\sigma}^2$

③ Quantile-quantile plot can assess normality, plots theoretical quantiles of a normal against empirical quantiles of (standardized) $\hat{\varepsilon}$.

Empirical quant of ε̂_i

theor. quant. of Norm

→ identity line

If $\{\varepsilon_i\}$ are approx normal

Eq uant

1.96 theor N quant

Heavy-tailed residuals