а) Как было показано на лекции,

$$\hat{\theta} = (X^T X + \lambda I_d)^{-1} X^T y$$

Тогда $MSE_{\hat{\theta}}(\theta) = \mathbb{E}[(\hat{\theta}-\theta)^T(\hat{\theta}-\theta)] =$

$$= \mathbb{E}[\hat{\theta}^T\hat{\theta} - 2\hat{\theta}^T\theta + \theta^T\theta] =$$

$$= \mathbb{E}\Big[ y^T X \underbrace{(X^T X + \lambda I_d)^{-1}(X^T X + \lambda I_d)^{-1}}_{W_\lambda} X^T y -$$

$$- 2\big((X^T X + \lambda I_d)^{-1} X^T y\big)^T \theta + \theta^T\theta\Big] =$$

$$= \mathbb{E}[y^T X W_\lambda^2 X^T y - 2\theta^T W_\lambda X^T y + \theta^T\theta]$$

$$= \mathbb{E}[tr(X W_\lambda^2 X^T \cdot y\, y^T) - 2\theta^T W_\lambda X^T y + \theta^T\theta] =$$

$$= tr(X W_\lambda^2 X^T \cdot \mathbb{E}[y y^T]) - 2\theta^T W_\lambda X^T \mathbb{E}[y] + \theta^T\theta =$$

$$= \{ \text{в предположении, что модель гауссовская} \}$$

$$= tr(X W_\lambda^2 X^T \sigma^2) - 2\theta^T W_\lambda X^T X\theta + \theta^T\theta$$

$$= tr\big(X W_\lambda^2 X^T \cdot (Cov\, y + (\mathbb{E}y)^T \mathbb{E}y)\big) - 2\theta^T W_\lambda X^T X\theta$$

$$+ \theta^T\theta = tr\big(X W_\lambda^2 X^T \cdot \sigma^2 + (X\theta)^T \cdot X\theta\big) - 2\theta^T W_\lambda X^T X\theta$$

$$+ \theta^T\theta$$

$$\ominus \{ \text{а также т.к. } Cov\, y = \mathbb{E}[(y - \mathbb{E}y)(y - \mathbb{E}y)^T] =$$

$$= \mathbb{E}[y y^T - 2y(\mathbb{E}y)^T + \mathbb{E}y \cdot (\mathbb{E}y)^T] \} \ominus$$

$$\ominus \, tr\big(X W_\lambda^2 X^T \cdot (Cov\, y + 2\mathbb{E}y(\mathbb{E}y)^T - \mathbb{E}y(\mathbb{E}y)^T)\big) -$$

$$- 2\theta^T W_\lambda X^T \mathbb{E}y + \theta^T\theta = tr\big(X W_\lambda^2 X^T \cdot (\sigma^2 I_n +$$

$$+ X\theta \cdot (X\theta)^T)) - 2\theta^T W_2 X^T \cdot X\theta + \theta^T\theta$$

Проверка: при $\lambda = 0$ имеем $W_2 = (X^T X)^{-1} =$

$$\Rightarrow MSE_{\hat\theta}(\theta) = tz\left( X(X^T X)^{-1} X^T \cdot (\sigma^2 I_n + X\theta (X\theta)^T) \right)$$

$$- 2\theta^T\theta + \theta^T\theta = tz\left( X(X^T X)^{-1} X^T \sigma^2 + X(X^T X)^{-1} X^T X\theta \cdot \right.$$

$$\cdot \theta^T X^T \right) - \theta^T\theta = tz\left( X(X^T X)^{-1}\sigma^2 + X\theta \cdot \theta^T X^T \right)$$

$$- \theta^T\theta$$

б) $<\hat{y}, e> = <\hat{y}, y - \hat{y}> = \left( X \cdot \underbrace{(X^T X + \lambda I_d)^{-1} X^T y}_{W_2} \right)^T \cdot y -$

$$- \left( X(X^T X + \lambda I_d)^{-1} X^T y \right)^T \cdot (-//-) =$$

$$= y^T X W_2 X^T y - y^T X W_2 X^T \cdot (X W_2 X^T y) =$$

$$= y^T X W_2 \underbrace{(I_d - X^T X W_2)}_{\neq 0} X^T y \neq 0, \quad z.\varepsilon.g.$$

θ Сначала одномерный случай:

θ $Pr_{\lambda x^2}(p) = \underset{x}{\operatorname{argmin}} \left( \underbrace{\frac{1}{2}(x-p)^2 + \lambda x^2}_{g(x,p)} \right)$

$\partial g(x,p) = x - p + 2\lambda x$

Т.к. $x_0 = Pr_{\lambda x^2}(p) \Leftrightarrow$ ~~0 ∈ p - x_0 ∈~~ $0 \in \partial g(x_0, p)$,

имеем: $0 \in \partial g(x_0, p) \Leftrightarrow 0 = x_0 - p + 2\lambda x_0 \Leftrightarrow$

$\Leftrightarrow x_0 = \dfrac{p}{1+2\lambda}$ — всегда ∃-т, т.к. $\lambda \geq 0$.

$\Rightarrow Pr_{\lambda x^2}(p) = \dfrac{p}{1+2\lambda}$.

По №6-ю 5 с лекции: если $\theta$ -реш-е

$F(\theta) + \lambda R(\theta) \to \min\limits_{\theta}$, то $\theta = P_{z_{\lambda R}}(\theta - \nabla F(\theta))$

В предположении, что есть сх-ть итер. метода, имеем: $\theta_{k+1} = \dfrac{\theta_k - \nabla F(\theta_k)}{1 + 2\lambda} =$

$= \dfrac{\theta_k - \nabla \|y - X\overline{\theta_k}\|_2^2}{1 + 2\lambda}$, где $\overline{\theta_k} = \begin{pmatrix} 0 \\ \vdots \\ \theta_k \\ \vdots \\ 0 \end{pmatrix} \overset{k-я}{\underset{\text{столбец}}{}} \leftarrow k\text{-я позиция}$

$\theta_{k+1} = \dfrac{\theta_k + 2 \, y^T_0 [X]_k - 2\theta_k}{1 + 2\lambda} = \dfrac{2 \, y^T [X]_k - \theta_k}{1 + 2\lambda}$

Многом. случай: $P_{z_{\lambda\|x\|_2^2}}(p) = \arg\min\limits_{x} \left( \frac{1}{2} \right.$

$\cdot \|x - p\|_2^2 + \lambda \|x\|_2^2 \big) = \arg\min\limits_{x} \left( \frac{1}{2} \sum\limits_{i=1}^{d} (x_i - p_i)^2 + \right.$

$+ \lambda \sum\limits_{i=1}^{d} x_i^2 = \left( \ldots, P_{z_{\lambda x_i^2}}(p_i), \ldots \right)^T$

Св-ва $\lambda$: 1) отличие от лассо-регрессии в том, что её итерации "плавные", т.к. нет резкого зануления $\theta_k$; а ещё сх-ть более быстрая

2) Сх-ть геометрическая, т.к. в знаменателе каждый раз $(1 + 2\lambda)$ стоит $\Rightarrow$ при $\lambda \to \infty$ очень быстро $\theta_k \to 0$, а при $\lambda = 0$ получается градиентный спуск для $F(\theta)$.