

On The Direct Maximization of Quadratic Weighted Kappa

David Vaughn
Derek Justice

DVAUGHN@MEASINC.COM
DJJUSTICE@MEASINC.COM

Abstract

In recent years, quadratic weighted kappa has been growing in popularity in the machine learning community as an evaluation metric in domains where the target labels to be predicted are drawn from integer ratings, usually obtained from human experts. For example, it was the metric of choice in several recent, high profile machine learning contests hosted on Kaggle :

www.kaggle.com/c/asap-aes ,
www.kaggle.com/c/asap-sas ,
www.kaggle.com/c/diabetic-retinopathy-detection .

Yet, little is understood about the nature of this metric, its underlying mathematical properties, where it fits among other common evaluation metrics such as mean squared error (MSE) and correlation, or if it can be optimized analytically, and if so, how. Much of this is due to the cumbersome way that this metric is commonly defined.

In this paper we first derive an equivalent but much simpler, and more useful, definition for quadratic weighted kappa, and then employ this alternate form to address the above issues.

1. Preliminaries

Although first developed in the statistical community as a measure of *inter-rater agreement*, κ has more recently become a popular performance metric in supervised machine learning, specifically in situations where the target (dependent) variable y is a discrete, interval variable (usually drawn from non-negative integers) such as is common in most human rating scales (e.g. “on a scale from 1 to 10”).

This differs from the ordinal regression setting, where there only exists an ordering over labels, but no intrinsic or constant length interval between them. Some have argued that the use of quadratic weighted kappa as a metric in the domain of human ratings imposes the erroneous assumption of “equal intervals” where there should be none (how this assumption is expressed in the metric itself will be made

clear in the following section). For example, when rating student essays on a scale from 1 to 5, the difference between a 1 and a 2 may not be equal to the difference between a 4 and a 5. While this may or may not be true in certain cases, we will not be concerned with that here.

1.1. Standard Definition

Quadratic weighted kappa, which we write κ to distinguish from linear weighted kappa, was originally developed as a measure of *inter-rater agreement*. In this scenario, there are two raters, \mathcal{A} and \mathcal{B} , each associated with a vector of n integer ratings $\mathbf{a}, \mathbf{b} \in \mathbb{L}^{n \times 1}$ where $\mathbb{L} = \{1, 2, \dots, \ell\}$ is a finite set of ℓ possible values. We seek to quantify the level of agreement between \mathbf{a} and \mathbf{b} . In order to compute $\kappa(\mathbf{a}, \mathbf{b})$, it is customary to start by computing frequency tables. The *observed* confusion matrix $\mathbf{U} = (u_{i,j}) \in \mathbb{N}^{\ell \times \ell}$ is first computed as:

$$u_{i,j} = \sum_{k=1}^n \mathbb{I}(a_k = i) \cdot \mathbb{I}(b_k = j)$$

Next, the *expected* confusion matrix $\mathbf{V} = (v_{i,j}) \in \mathbb{R}^{\ell \times \ell}$ is computed by assuming there is no correlation between raters. Under this assumption, we can simply compute \mathbf{V} as the outer product between the two rater’s observed marginal distributions, normalized so that \mathbf{V} has the same sum as \mathbf{U} (i.e. $\sum_{i,j=1}^{\ell} v_{i,j} = \sum_{i,j=1}^{\ell} u_{i,j}$):

$$\mathbf{V} = \frac{(\mathbf{U} \cdot \mathbf{e}) \otimes (\mathbf{U}^T \cdot \mathbf{e})}{n}$$

where \mathbf{e} denotes the all ones vector. Finally, a matrix of weights $\mathbf{W} = (w_{i,j}) \in \mathbb{R}^{\ell \times \ell}$ is defined as:

$$w_{i,j} = \frac{(i - j)^2}{(\ell - 1)^2}$$

Given \mathbf{U} , \mathbf{V} , and \mathbf{W} , it is customary to define κ as:

$$\kappa(\mathbf{a}, \mathbf{b}) = 1 - \frac{\sum_{i,j=1}^{\ell} u_{i,j} \cdot w_{i,j}}{\sum_{i,j=1}^{\ell} v_{i,j} \cdot w_{i,j}} = 1 - \frac{\langle \mathbf{U}, \mathbf{W} \rangle_F}{\langle \mathbf{V}, \mathbf{W} \rangle_F} \quad (1)$$

where $\langle \cdot, \cdot \rangle_F$ denotes the Frobenius inner product.

1.2. Alternate Form

Assume a data matrix $\mathbf{X} = [\mathbf{x}_1^\top, \dots, \mathbf{x}_n^\top]^\top \in \mathbb{R}^{n \times d}$ of n points in d dimensional space, and a vector $\mathbf{y} \in \mathbb{L}^{n \times 1}$ of n integer labels where $\mathbb{L} = \{1, 2, \dots, \ell\}$ is a finite set of possible labels. We assume there is some true functional relationship $y_i = f(x_i)$ between each point x_i and it's label y_i . The goal is to find a function $\hat{y}_i = \hat{f}(x_i)$ which approximates the true function as closely as possible. Using this notation, we can begin to rewrite the standard definition of κ by first noting that the numerator simply represents a standard sum of squared errors:

$$\langle \mathbf{U}, \mathbf{W} \rangle_F = \sum_{k=1}^n \frac{(y_k - \hat{f}(x_k))^2}{(\ell - 1)^2} = \frac{\|\mathbf{y} - \hat{\mathbf{y}}\|^2}{(\ell - 1)^2}$$

Next, we note that the denominator can be similarly rewritten in terms of \mathbf{y} and $\hat{\mathbf{y}}$:

$$\langle \mathbf{V}, \mathbf{W} \rangle_F = \frac{\|\mathbf{y}\|^2 - \frac{2}{n}(\mathbf{y}^\top \mathbf{e})(\hat{\mathbf{y}}^\top \mathbf{e}) + \|\hat{\mathbf{y}}\|^2}{(\ell - 1)^2}$$

By combining these expressions we can derive a simplified expression for κ purely in terms of \mathbf{y} and $\hat{\mathbf{y}}$, without the need for contingency tables:

$$\begin{aligned} \kappa(\mathbf{y}, \hat{\mathbf{y}}) &= 1 - \frac{\|\mathbf{y} - \hat{\mathbf{y}}\|^2}{\|\mathbf{y}\|^2 - \frac{2}{n}(\mathbf{y}^\top \mathbf{e})(\hat{\mathbf{y}}^\top \mathbf{e}) + \|\hat{\mathbf{y}}\|^2} \\ &= 1 - \frac{\|\mathbf{y}\|^2 - 2\mathbf{y}^\top \hat{\mathbf{y}} + \|\hat{\mathbf{y}}\|^2}{\|\mathbf{y}\|^2 - \frac{2}{n}(\mathbf{y}^\top \mathbf{e})(\hat{\mathbf{y}}^\top \mathbf{e}) + \|\hat{\mathbf{y}}\|^2} \\ &= \frac{2[\mathbf{y} - \mathbf{e}(\frac{\mathbf{y}^\top \mathbf{e}}{n})]^\top \hat{\mathbf{y}}}{\|\mathbf{y}\|^2 - \frac{2}{n}(\mathbf{y}^\top \mathbf{e})(\hat{\mathbf{y}}^\top \mathbf{e}) + \|\hat{\mathbf{y}}\|^2} \end{aligned}$$

If we drop the integer constraint on labels then we can assume that \mathbf{y} has been centered, which allows further simplification:

$$\kappa(\mathbf{y}, \hat{\mathbf{y}}) = \frac{2\langle \mathbf{y}, \hat{\mathbf{y}} \rangle}{\|\mathbf{y}\|^2 + \|\hat{\mathbf{y}}\|^2} \quad (2)$$

Throughout the rest of this paper we will assume (without loss of generality) that \mathbf{y} is centered.

2. Linear Model

We consider a simple linear model for $\mathbf{y} = f(\mathbf{X})$:

$$f(\mathbf{X}) = \mathbf{X}\alpha$$

and so we seek the function parameter $\hat{\alpha}$ that maximizes $\kappa(\alpha)$, where

$$\kappa(\alpha) = \frac{2\langle \mathbf{y}, \mathbf{X}\alpha \rangle}{\|\mathbf{y}\|^2 + \|\mathbf{X}\alpha\|^2}$$

We will rely on the fact that $\kappa(\alpha)$ is *strictly quasiconcave* (see appendix for proof), since it is known that any local maximum of a *strictly quasiconcave* function f is also a global maximum (Avriel, 2003). And so we can apply the method of Lagrange multipliers to find a local maximum, and conclude that this is the global maximum as well. To begin, we can use the variable substitution $\gamma = \|\mathbf{y}\|^2 + \|\mathbf{X}\alpha\|^2$ to write this as a constrained optimization problem:

$$\begin{aligned} \max_{\alpha, \gamma} \quad & \frac{2}{\gamma} \langle \mathbf{y}, \mathbf{X}\alpha \rangle \\ \text{subject to} \quad & \gamma = \|\mathbf{y}\|^2 + \|\mathbf{X}\alpha\|^2 \end{aligned}$$

From here we can define the Lagrangian:

$$\begin{aligned} \mathbf{L}(\alpha, \gamma, \lambda) &= \frac{2}{\gamma} \langle \mathbf{y}, \mathbf{X}\alpha \rangle + \lambda(\gamma - \|\mathbf{y}\|^2 - \|\mathbf{X}\alpha\|^2) \\ &= \frac{2}{\gamma} \mathbf{y}^\top \mathbf{X}\alpha + \lambda(\gamma - \mathbf{y}^\top \mathbf{y} - \alpha^\top \mathbf{X}^\top \mathbf{X}\alpha) \end{aligned}$$

We solve by first differentiating w.r.t. γ and setting equal to zero:

$$\begin{aligned} \frac{\partial \mathbf{L}}{\partial \gamma} \Big|_{\hat{\alpha}, \hat{\gamma}, \hat{\lambda}} &= -\frac{2}{\hat{\gamma}^2} \mathbf{y}^\top \mathbf{X}\hat{\alpha} + \hat{\lambda} = 0 \\ \hat{\lambda} &= \frac{2}{\hat{\gamma}^2} \mathbf{y}^\top \mathbf{X}\hat{\alpha} \quad (3) \end{aligned}$$

Then w.r.t. α :

$$\frac{\partial \mathbf{L}}{\partial \alpha} \Big|_{\hat{\alpha}, \hat{\gamma}, \hat{\lambda}} = \frac{2}{\hat{\gamma}} \mathbf{X}^\top \mathbf{y} - 2\hat{\lambda} \mathbf{X}^\top \mathbf{X}\hat{\alpha} = 0$$

Substituting in the previous expression (3) for $\hat{\lambda}$ yields:

$$\begin{aligned} \frac{\mathbf{X}^\top \mathbf{y}}{\hat{\gamma}} - \frac{2\mathbf{y}^\top \mathbf{X}\hat{\alpha}}{\hat{\gamma}^2} \mathbf{X}^\top \mathbf{X}\hat{\alpha} &= 0 \\ \frac{2\mathbf{y}^\top \mathbf{X}\hat{\alpha}}{\hat{\gamma}} \mathbf{X}^\top \mathbf{X}\hat{\alpha} &= \mathbf{X}^\top \mathbf{y} \end{aligned}$$

Now substituting back in for $\hat{\gamma}$ gives a surprising result:

$$\begin{aligned} \frac{2\mathbf{y}^\top \mathbf{X}\hat{\alpha}}{\|\mathbf{y}\|^2 + \|\mathbf{X}\hat{\alpha}\|^2} \mathbf{X}^\top \mathbf{X}\hat{\alpha} &= \mathbf{X}^\top \mathbf{y} \\ \hat{\kappa} \mathbf{X}^\top \mathbf{X}\hat{\alpha} &= \mathbf{X}^\top \mathbf{y} \\ \hat{\alpha} &= \frac{1}{\hat{\kappa}} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad (4) \end{aligned}$$

where we have simplified by recognizing the expression for $\hat{\kappa} = \kappa(\hat{\alpha})$. In addition, we recognize the “least-squares” solution $\hat{\alpha}_{\ell s} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$, the well studied minimizer of the quantity $\|\mathbf{y} - \mathbf{X}\alpha\|$, *a.k.a.* the norm of the residual. This means that $\hat{\alpha}$ (which we will write $\hat{\alpha}_\kappa$ when necessary to avoid confusion) and the least-squares solution $\hat{\alpha}_{\ell s}$ are related via $\hat{\kappa}$ by a surprisingly simple formula:

$$\hat{\alpha}_\kappa = \frac{1}{\hat{\kappa}} \hat{\alpha}_{\ell s} \quad (5)$$

In order to obtain a closed form solution for $\hat{\alpha}_\kappa$ we must solve for $\hat{\kappa}$ solely in terms of \mathbf{X} and \mathbf{y} . We employ another useful tool from least-squares regression, the “hat” matrix $\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$. The hat matrix allows the least-squares labels $\hat{\mathbf{y}}_{\ell s}$ to be expressed solely in terms of the original labels \mathbf{y} :

$$\hat{\mathbf{y}}_{\ell s} = \mathbf{X}\hat{\alpha}_{\ell s} = \mathbf{H}\mathbf{y} \quad (6)$$

Using the well-known fact that the hat matrix is both symmetric and idempotent ($\mathbf{H}^\top \mathbf{H} = \mathbf{H}^2 = \mathbf{H}$), we can establish the following identity:

$$\begin{aligned} \|\hat{\mathbf{y}}_{\ell s}\|^2 &= (\mathbf{H}\mathbf{y})^\top (\mathbf{H}\mathbf{y}) \\ &= \mathbf{y}^\top \mathbf{H}^\top \mathbf{H} \mathbf{y} \\ &= \mathbf{y}^\top \mathbf{H} \mathbf{y} \\ &= \mathbf{y}^\top \mathbf{X} \hat{\alpha}_{\ell s} = \langle \mathbf{y} \cdot \mathbf{X} \hat{\alpha}_{\ell s} \rangle \end{aligned} \quad (7)$$

Now we can take the derived expression for $\hat{\alpha}_\kappa$ from (5) and substitute it back into the original expression for κ and apply the above identity (7) to solve for $\hat{\kappa}$:

$$\begin{aligned} \hat{\kappa} = \kappa(\hat{\alpha}_\kappa) &= \frac{2\langle \mathbf{y} \cdot \mathbf{X} \hat{\alpha}_\kappa \rangle}{\|\mathbf{y}\|^2 + \|\mathbf{X} \hat{\alpha}_\kappa\|^2} \\ &= \frac{2\langle \mathbf{y} \cdot \frac{1}{\hat{\kappa}} \mathbf{X} \hat{\alpha}_{\ell s} \rangle}{\|\mathbf{y}\|^2 + \|\frac{1}{\hat{\kappa}} \mathbf{X} \hat{\alpha}_{\ell s}\|^2} \\ &= \frac{\frac{2}{\hat{\kappa}} \langle \mathbf{y} \cdot \mathbf{X} \hat{\alpha}_{\ell s} \rangle}{\|\mathbf{y}\|^2 + \frac{1}{\hat{\kappa}^2} \|\mathbf{X} \hat{\alpha}_{\ell s}\|^2} \\ &= \frac{2\hat{\kappa} \|\hat{\mathbf{y}}_{\ell s}\|^2}{\hat{\kappa}^2 \|\mathbf{y}\|^2 + \|\hat{\mathbf{y}}_{\ell s}\|^2} \\ \hat{\kappa}^2 \|\mathbf{y}\|^2 + \|\hat{\mathbf{y}}_{\ell s}\|^2 &= 2\|\hat{\mathbf{y}}_{\ell s}\|^2 \\ \hat{\kappa} &= \frac{\|\hat{\mathbf{y}}_{\ell s}\|}{\|\mathbf{y}\|} \end{aligned} \quad (8)$$

Finally we can express $\hat{\alpha}$ solely in terms of \mathbf{X} and \mathbf{y} :

$$\begin{aligned} \hat{\alpha}_\kappa &= \frac{1}{\hat{\kappa}} \hat{\alpha}_{\ell s} = \frac{\|\mathbf{y}\|}{\|\hat{\mathbf{y}}_{\ell s}\|} \hat{\alpha}_{\ell s} = \frac{\|\mathbf{y}\|}{\|\hat{\mathbf{y}}_{\ell s}\|} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\ &= \left(\sqrt{\frac{\mathbf{y}^\top \mathbf{y}}{\mathbf{y}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}}} \right) (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\ &= \left(\sqrt{\frac{\mathbf{y}^\top \mathbf{y}}{\mathbf{y}^\top \mathbf{H} \mathbf{y}}} \right) (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \end{aligned} \quad (9)$$

2.1. Regularization

It is often useful to include a regularization term in the optimization to account for cases where $\mathbf{X}^\top \mathbf{X}$ is singular or ill-conditioned, which could lead to a numerically unstable (or non-existent) solution. In addition, it provides a convenient mechanism for manually tuning the bias/variance trade-off. It is easy to accommodate a norm penalty term within the current framework, yielding a slightly modified optimization problem:

$$\begin{aligned} \max_{\alpha, \gamma} \quad & \frac{2}{\gamma} \langle \mathbf{y} \cdot \mathbf{X} \alpha \rangle - \mu \|\alpha\|^2 \\ \text{subject to} \quad & \gamma = \|\mathbf{y}\|^2 + \|\mathbf{X} \alpha\|^2 \end{aligned}$$

where μ controls the magnitude of the penalty, and is usually set experimentally (e.g. with cross-validation). This leads to a similarly modified Lagrangian:

$$\mathbf{L}(\alpha, \gamma, \lambda) = \frac{2}{\gamma} \mathbf{y}^\top \mathbf{X} \alpha - \mu \alpha^\top \alpha + \lambda (\gamma - \mathbf{y}^\top \mathbf{y} - \alpha^\top \mathbf{X}^\top \mathbf{X} \alpha)$$

Proceeding in the same manner as the last section, differentiating w.r.t. γ and setting equal to zero yields the same result (3) for $\hat{\lambda}$. Differentiating w.r.t. α yields a slightly different result:

$$\left. \frac{\partial \mathbf{L}}{\partial \alpha} \right|_{\hat{\alpha}, \hat{\gamma}, \hat{\lambda}} = \frac{2}{\hat{\gamma}} \mathbf{X}^\top \mathbf{y} - 2\hat{\lambda} \mathbf{X}^\top \mathbf{X} \hat{\alpha} - 2\mu \hat{\alpha} = 0$$

Again, substituting in for $\hat{\lambda}$ and manipulating the result, we end up with:

$$\begin{aligned} \hat{\kappa} \mathbf{X}^\top \mathbf{X} \hat{\alpha} + \mu \hat{\gamma} \hat{\alpha} &= \mathbf{X}^\top \mathbf{y} \\ \hat{\alpha} &= \frac{1}{\hat{\kappa}} (\mathbf{X}^\top \mathbf{X} + \frac{\mu \hat{\gamma}}{\hat{\kappa}} \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y} \end{aligned}$$

In order to obtain a closed form solution for $\hat{\alpha}$ all that is important to recognize is that we once again know the optimal solution to within a factor $\varphi = \frac{1}{\hat{\kappa}}$, just as in the regularized case. The only difference is that now, instead of $\hat{\alpha}$ just being within a scalar multiple of the least-squares solution,

$\hat{\alpha}$ is now within a scalar multiple of the *regularized* least squares solution, i.e. the *ridge regression* solution $\hat{\alpha}_{rr}$. As μ is just a regularization tuning parameter, it can simply absorb the $\hat{\kappa}$ and $\hat{\gamma}$. We make the substitution $\sigma = \frac{\mu\hat{\gamma}}{\hat{\kappa}}$ to make this more clear:

$$\hat{\alpha}_{\kappa} = \frac{1}{\hat{\kappa}}(\mathbf{X}^T \mathbf{X} + \sigma \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} = \frac{1}{\hat{\kappa}} \hat{\alpha}_{rr} \quad (10)$$

Just as in the previous case, we can take the derived expression for $\hat{\alpha}_{\kappa}$ from (10) and substitute it back into the original expression for κ to solve for $\hat{\kappa}$:

$$\begin{aligned} \hat{\kappa} = \kappa(\hat{\alpha}_{\kappa}) &= \frac{2\langle \mathbf{y} \cdot \mathbf{X} \hat{\alpha}_{\kappa} \rangle}{\|\mathbf{y}\|^2 + \|\mathbf{X} \hat{\alpha}_{\kappa}\|^2} \\ &= \frac{2\langle \mathbf{y} \cdot \frac{1}{\hat{\kappa}} \mathbf{X} \hat{\alpha}_{rr} \rangle}{\|\mathbf{y}\|^2 + \|\frac{1}{\hat{\kappa}} \mathbf{X} \hat{\alpha}_{rr}\|^2} \\ &= \frac{\frac{2}{\hat{\kappa}} \langle \mathbf{y} \cdot \mathbf{X} \hat{\alpha}_{rr} \rangle}{\|\mathbf{y}\|^2 + \frac{1}{\hat{\kappa}^2} \|\mathbf{X} \hat{\alpha}_{rr}\|^2} \\ \hat{\kappa}^2 \|\mathbf{y}\|^2 + \|\mathbf{X} \hat{\alpha}_{rr}\|^2 &= 2\langle \mathbf{y} \cdot \mathbf{X} \hat{\alpha}_{rr} \rangle \\ \hat{\kappa} &= \frac{\sqrt{2\langle \mathbf{y} \cdot \mathbf{X} \hat{\alpha}_{rr} \rangle - \|\mathbf{X} \hat{\alpha}_{rr}\|^2}}{\|\mathbf{y}\|} \quad (11) \end{aligned}$$

In this case the result does not simplify as nicely as in the non-regularized case, due to the fact that \mathbf{H}_{σ} (where $\mathbf{H}_{\sigma} = \mathbf{X}(\mathbf{X}^T \mathbf{X} + \sigma \mathbf{I})^{-1} \mathbf{X}^T$ denotes the regularized variant of the “hat” matrix) is no longer idempotent. Still, we can use \mathbf{H}_{σ} to write a somewhat simplified expression for $\hat{\kappa}$:

$$\hat{\kappa} = \sqrt{\frac{2\mathbf{y}^T \mathbf{H}_{\sigma} \mathbf{y} - \mathbf{y}^T \mathbf{H}_{\sigma}^2 \mathbf{y}}{\mathbf{y}^T \mathbf{y}}} \quad (12)$$

Finally we can express the regularized solution $\hat{\alpha}_{\kappa}$ in closed form:

$$\begin{aligned} \hat{\alpha}_{\kappa} &= \frac{1}{\hat{\kappa}} \hat{\alpha}_{rr} \\ &= \left(\sqrt{\frac{\mathbf{y}^T \mathbf{y}}{2\mathbf{y}^T \mathbf{H}_{\sigma} \mathbf{y} - \mathbf{y}^T \mathbf{H}_{\sigma}^2 \mathbf{y}}} \right) (\mathbf{X}^T \mathbf{X} + \sigma \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \quad (13) \end{aligned}$$

3. Link to Correlation

There is an interesting link between quadratic weighted kappa and correlation. There are many (equivalent) definitions of correlation, but the simplest is just the normalized dot product between two vectors \mathbf{u}, \mathbf{v} :

$$\rho(\mathbf{u}, \mathbf{v}) = \frac{\langle \mathbf{u} \cdot \mathbf{v} \rangle}{\|\mathbf{u}\| \|\mathbf{v}\|}$$

Now returning to the (non-regularized) κ -optimal solution labels $\hat{\mathbf{y}}_{\kappa} = \mathbf{X} \hat{\alpha}_{\kappa}$, we can expand the expression for $\rho(\mathbf{y}, \hat{\mathbf{y}}_{\kappa})$ and apply (7) and (8) to simplify:

$$\begin{aligned} \rho(\mathbf{y}, \hat{\mathbf{y}}_{\kappa}) &= \frac{\langle \mathbf{y} \cdot \hat{\mathbf{y}}_{\kappa} \rangle}{\|\mathbf{y}\| \|\hat{\mathbf{y}}_{\kappa}\|} \\ &= \frac{\langle \mathbf{y} \cdot \mathbf{X} \hat{\alpha}_{\kappa} \rangle}{\|\mathbf{y}\| \|\mathbf{X} \hat{\alpha}_{\kappa}\|} \\ &= \frac{\frac{1}{\hat{\kappa}} \langle \mathbf{y} \cdot \mathbf{X} \hat{\alpha}_{rr} \rangle}{\|\mathbf{y}\| \|\frac{1}{\hat{\kappa}} \mathbf{X} \hat{\alpha}_{rr}\|} \\ &= \frac{\frac{1}{\hat{\kappa}} \|\hat{\mathbf{y}}_{\kappa}\|^2}{\frac{1}{\hat{\kappa}} \|\mathbf{y}\| \|\hat{\mathbf{y}}_{\kappa}\|} \\ &= \frac{\|\hat{\mathbf{y}}_{\kappa}\|}{\|\mathbf{y}\|} = \kappa(\mathbf{y}, \hat{\mathbf{y}}_{\kappa}) = \hat{\kappa} \quad (14) \end{aligned}$$

And so while correlation is, in general, not equivalent to quadratic weighted kappa, they do coincide precisely at the point where kappa is maximized.

Appendix: Proof of Strict Quasiconcavity

We wish to show that the function

$$\kappa(\alpha) = \frac{2\langle \mathbf{y} \cdot \mathbf{X}\alpha \rangle}{\|\mathbf{y}\|^2 + \|\mathbf{X}\alpha\|^2}$$

is *strictly quasiconcave*. A function f is said to be *strictly quasiconcave* if, for any two points \mathbf{u}, \mathbf{v} and any positive weight $\theta < 1$, the following is true (Avriel, 2003):

$$f(\theta\mathbf{u} + (1 - \theta)\mathbf{v}) > \min[f(\mathbf{u}), f(\mathbf{v})] \quad (\text{A-1})$$

Let us assume, without loss of generality, that $\min[\kappa(\mathbf{u}), \kappa(\mathbf{v})] = \kappa(\mathbf{u})$, which leads us to write:

$$\begin{aligned} \kappa(\mathbf{v}) &\geq \kappa(\mathbf{u}) \\ \frac{2\langle \mathbf{y} \cdot \mathbf{X}\mathbf{v} \rangle}{\|\mathbf{y}\|^2 + \|\mathbf{X}\mathbf{v}\|^2} &\geq \frac{2\langle \mathbf{y} \cdot \mathbf{X}\mathbf{u} \rangle}{\|\mathbf{y}\|^2 + \|\mathbf{X}\mathbf{u}\|^2} \\ \mathbf{y}^\top \mathbf{X}\mathbf{v} &\geq \mathbf{y}^\top \mathbf{X}\mathbf{u} \left(\frac{\|\mathbf{y}\|^2 + \|\mathbf{X}\mathbf{v}\|^2}{\|\mathbf{y}\|^2 + \|\mathbf{X}\mathbf{u}\|^2} \right) \end{aligned} \quad (\text{A-2})$$

Returning to the the definition of *strictly quasiconcave* and expanding the expression for $\kappa(\theta\mathbf{u} + (1 - \theta)\mathbf{v})$ gives us:

$$\kappa(\theta\mathbf{u} + (1 - \theta)\mathbf{v}) = \frac{2\theta\mathbf{y}^\top \mathbf{X}\mathbf{u} + 2(1 - \theta)\mathbf{y}^\top \mathbf{X}\mathbf{v}}{\|\mathbf{y}\|^2 + \|\theta\mathbf{X}\mathbf{u} + (1 - \theta)\mathbf{X}\mathbf{v}\|^2}$$

Now, we know that the function $h(\mathbf{z}) = \|\mathbf{X}\mathbf{z}\|^2$ is strictly convex (because $\|\mathbf{X}\mathbf{z}\|^2 = \mathbf{z}^\top \mathbf{X}^\top \mathbf{X}\mathbf{z}$, where $\mathbf{X}^\top \mathbf{X} \succ 0$ since $\mathbf{X}^\top \mathbf{X}$ is assumed non-singular). By the definition of strict convexity, this means:

$$\|\theta\mathbf{X}\mathbf{u} + (1 - \theta)\mathbf{X}\mathbf{v}\|^2 < \theta\|\mathbf{X}\mathbf{u}\|^2 + (1 - \theta)\|\mathbf{X}\mathbf{v}\|^2$$

Therefore we can substitute the expression on the right into the previous equation to obtain:

$$\kappa(\theta\mathbf{u} + (1 - \theta)\mathbf{v}) > \frac{2\theta\mathbf{y}^\top \mathbf{X}\mathbf{u} + 2(1 - \theta)\mathbf{y}^\top \mathbf{X}\mathbf{v}}{\|\mathbf{y}\|^2 + \theta\|\mathbf{X}\mathbf{u}\|^2 + (1 - \theta)\|\mathbf{X}\mathbf{v}\|^2} \quad (\text{A-3})$$

We can then substitute the expression for $\mathbf{y}^\top \mathbf{X}\mathbf{v}$ from the right side of (A-2) into the above, leading to:

$$\begin{aligned} \frac{2\theta\mathbf{y}^\top \mathbf{X}\mathbf{u} + 2(1 - \theta)\mathbf{y}^\top \mathbf{X}\mathbf{v}}{\|\mathbf{y}\|^2 + \theta\|\mathbf{X}\mathbf{u}\|^2 + (1 - \theta)\|\mathbf{X}\mathbf{v}\|^2} &\geq \\ \frac{2\theta\mathbf{y}^\top \mathbf{X}\mathbf{u} + 2(1 - \theta)\mathbf{y}^\top \mathbf{X}\mathbf{u} \left(\frac{\|\mathbf{y}\|^2 + \|\mathbf{X}\mathbf{v}\|^2}{\|\mathbf{y}\|^2 + \|\mathbf{X}\mathbf{u}\|^2} \right)}{\|\mathbf{y}\|^2 + \theta\|\mathbf{X}\mathbf{u}\|^2 + (1 - \theta)\|\mathbf{X}\mathbf{v}\|^2} &= \\ \frac{2\mathbf{y}^\top \mathbf{X}\mathbf{u} \left(\theta + (1 - \theta) \left(\frac{\|\mathbf{y}\|^2 + \|\mathbf{X}\mathbf{v}\|^2}{\|\mathbf{y}\|^2 + \|\mathbf{X}\mathbf{u}\|^2} \right) \right)}{\|\mathbf{y}\|^2 + \theta\|\mathbf{X}\mathbf{u}\|^2 + (1 - \theta)\|\mathbf{X}\mathbf{v}\|^2} &= \\ \frac{2\mathbf{y}^\top \mathbf{X}\mathbf{u} \left(\frac{\theta(\|\mathbf{y}\|^2 + \|\mathbf{X}\mathbf{u}\|^2) + (1 - \theta)(\|\mathbf{y}\|^2 + \|\mathbf{X}\mathbf{v}\|^2)}{\|\mathbf{y}\|^2 + \|\mathbf{X}\mathbf{u}\|^2} \right)}{\|\mathbf{y}\|^2 + \theta\|\mathbf{X}\mathbf{u}\|^2 + (1 - \theta)\|\mathbf{X}\mathbf{v}\|^2} &= \\ \frac{\frac{2\mathbf{y}^\top \mathbf{X}\mathbf{u}}{\|\mathbf{y}\|^2 + \|\mathbf{X}\mathbf{u}\|^2} (\|\mathbf{y}\|^2 + \theta\|\mathbf{X}\mathbf{u}\|^2 + (1 - \theta)\|\mathbf{X}\mathbf{v}\|^2)}{\|\mathbf{y}\|^2 + \theta\|\mathbf{X}\mathbf{u}\|^2 + (1 - \theta)\|\mathbf{X}\mathbf{v}\|^2} &= \\ = \kappa(\mathbf{u}) \end{aligned}$$

Combining this result with (A-3) we have

$$\kappa(\theta\mathbf{u} + (1 - \theta)\mathbf{v}) > \kappa(\mathbf{u})$$

or, stated another way

$$\kappa(\theta\mathbf{u} + (1 - \theta)\mathbf{v}) > \min[\kappa(\mathbf{u}), \kappa(\mathbf{v})]$$

which is what we were trying to prove. \square

References

Avriel, Mordecai. *Nonlinear programming: analysis and methods*. Courier Corporation, 2003.